

Вероятностное тематическое моделирование

К. В. Воронцов

7 апреля 2013 г.

Содержание

1	Задачи тематического моделирования	2
1.1	Вероятностная модель коллекции документов	3
1.2	Предварительная обработка данных	5
1.3	Метод максимума правдоподобия	6
1.4	Внутренние оценки качества тематических моделей	7
2	Вероятностный латентный семантический анализ	9
2.1	EM-алгоритм	9
2.2	Обобщённый EM-алгоритм	11
2.3	Онлайновый EM-алгоритм	12
2.4	Стохастический EM-алгоритм	15
2.5	Формирование начальных приближений	16
2.6	Частичное обучение	17
3	Латентное размещение Дирихле	19
3.1	Байесовский вывод	19
3.2	Сэмплирование Гиббса	20
3.3	Оптимизация гиперпараметров	22
3.4	Позволяет ли сглаживание уменьшить переобучение	22
4	Робастные и разреженные тематические модели	23
4.1	Робастная тематическая модель с шумом и фоном	23
4.2	Принудительное разреживание	27
5	Иерархические тематические модели	29
5.1	Определение тематического дерева	30
5.2	Иерархическая модель для категоризации текстов	31
6	Критерии качества тематических моделей	34
6.1	Критерии условной независимости	34
6.2	Критерии качества классификации документов	37
6.3	Критерии качества тематического поиска	37

1 Задачи тематического моделирования

Тематическое моделирование (topic modeling) — одно из современных приложений машинного обучения к анализу текстов, активно развивающееся с конца 90-х годов. *Тематическая модель* (topic model) коллекции текстовых документов определяет, к каким темам относится каждый документ и какие слова (термины) образуют каждую тему.

Вероятностная тематическая модель (ВТМ) описывает каждую тему дискретным распределением на множестве терминов, каждый документ — дискретным распределением на множестве тем. Предполагается, что коллекция документов — это последовательность терминов, выбранных случайно и независимо из смеси таких распределений, и ставится задача восстановления компонент смеси по выборке. В большинстве приложений требуется определить также и число тем.

Поскольку документ или термин может относиться одновременно ко многим темам с различными вероятностями, говорят, что ВТМ осуществляет «мягкую» кластеризацию документов и терминов по кластерам-темам. Тем самым решаются проблемы синонимии и омонимии терминов, возникающие при обычной «жёсткой» кластеризации. Синонимы, часто употребляющиеся в схожих контекстах, с большой вероятностью попадают в одну тему. Омонимы, употребляющиеся в разных контекстах, распределяются между несколькими темами соответственно частоте употребления.

Тематические модели применяются для выявления трендов в научных публикациях или новостных потоках [44, 33], для классификации и категоризации документов [28] и изображений [20, 14], для семантического информационного поиска [41], в том числе многоязычного [34], для тегирования веб-страниц [19], для обнаружения текстового спама [4], в рекомендательных системах [40] и других приложениях. Для конкретности будем рассматривать применение ВТМ для категоризации и тематического поиска научных публикаций.

В информационном поиске документы принято представлять векторами, координаты которых соответствуют словам, а значения — статистическим характеристикам слов, например частотам или tf-idf. Поиск документов по коротким запросам реализуется путём поиска векторов, в которых часто встречаются слова запроса [3]. Тематическая модель позволяет использовать тот же механизм для поиска документов схожей тематики по целому документу или по длинному фрагменту текста. При этом документы представляются векторами тем, а не векторами слов. Векторами тем представляются также связанные с документами объекты: термины, рисунки, авторы, научные группы, организации, конференции, журналы, сайты и т. д., что позволяет задавать в качестве запроса любой объект или совокупность объектов и искать по ним объекты того же или другого типа, имеющие схожую тематику.

Тематические модели могут учитывать различные особенности языка и текстовых коллекций. Существуют модели, выявляющие устойчивые последовательности терминов, отслеживающие изменения тематики во времени или внутри отдельных документов, строящие иерархические отношения между темами, учитывающие связи между документами через авторство или ссылки, и т. д. Многочисленные разновидности вероятностных тематических моделей описаны в обзоре [11].

§1.1 Вероятностная модель коллекции документов

Пусть D — множество (коллекция) текстовых документов, W — множество (словарь) всех употребляемых в них терминов (слов или словосочетаний). Каждый документ $d \in D$ представляет собой последовательность n_d терминов (w_1, \dots, w_{n_d}) из словаря W . Термин может повторяться в документе много раз.

Вероятностное пространство и гипотеза независимости. Предполагается, что существует конечное множество тем T , и каждое употребление термина w в каждом документе d связано с некоторой темой $t \in T$, которая не известна. Коллекция документов рассматривается как множество троек (d, w, t) , выбранных *случайно и независимо* из дискретного распределения $p(d, w, t)$, заданного на конечном множестве $D \times W \times T$. Документы $d \in D$ и термины $w \in W$ являются наблюдаемыми переменными, тема $t \in T$ является *латентной* (скрытой) переменной.

Гипотеза о независимости элементов выборки эквивалентна предположению, что порядок терминов в документах не важен для выявления тематики, то есть тематику документа можно узнать даже после произвольной перестановки терминов, хотя для человека такой текст теряет смысл. Это предположение называют гипотезой «мешка слов» (bag of words). Порядок документов в коллекции также не имеет значения; это предположение называют гипотезой «мешка документов».

Приняв гипотезу «мешка слов», можно перейти к более компактному представлению документа как подмножества $d \subset W$, в котором каждому элементу $w \in d$ поставлено в соответствие число n_{dw} вхождений термина w в документ d .

Постановка задачи тематического моделирования. Построить *тематическую модель* коллекции документов D — значит найти множество тем T , распределения $p(w | t)$ для всех тем $t \in T$ и распределения $p(t | d)$ для всех документов $d \in D$. Можно также говорить о задаче совместной «мягкой» кластеризации множества документов и множества слов по множеству кластеров-тем. *Мягкая кластеризация* означает, что каждый документ или термин не жёстко приписывается какой-то одной теме, а распределяется по нескольким темам.

Найденные распределения используются затем для решения прикладных задач. Распределение $p(t | d)$ является удобным признаковым описанием документа в задачах информационного поиска, классификации и категоризации документов.

Гипотеза условной независимости. Будем полагать, что появление слов в документе d , относящихся к теме t , описывается общим для всей коллекции распределением $p(w | t)$ и не зависит от документа d . Это предположение, называемое *гипотезой условной независимости*, допускает три эквивалентных представления:

$$p(w | d, t) = p(w | t); \quad p(d | w, t) = p(d | t); \quad p(d, w | t) = p(d | t)p(w | t). \quad (1.1)$$

Вероятностная модель порождения данных. Согласно определению условной вероятности, формуле полной вероятности и гипотезе условной независимости

$$p(w | d) = \sum_{t \in T} p(t | d) p(w | t), \quad (1.2)$$

Алгоритм 1.1. Вероятностная модель порождения коллекции документов.

Вход: распределения $p(w | t)$, $p(t | d)$;

Выход: выборка пар (d_i, w_i) , $i = 1, \dots, n$;

- 1 **для всех** $d \in D$
 - 2 задать длину n_d документа d ;
 - 3 **для всех** $i = 1, \dots, n_d$
 - 4 выбрать случайную тему t из распределения $p(t | d)$;
 - 5 выбрать случайный термин w из распределения $p(w | t)$;
 - 6 добавить в выборку пару (d, w) , при этом тема t «забывается»;
-

где $p(t | d)$ и $p(w | t)$ — искомые распределения. Согласно порождающей модели (1.2), коллекция D — это выборка наблюдений (d, w) , генерируемых Алгоритмом 1.1.

Гипотеза разреженности. Естественно предполагать, что каждый документ d и каждый термин w связан с небольшим числом тем t . В таком случае значительная часть вероятностей $p(t | d)$ и $p(w | t)$ должна обращаться в нуль.

Если документ относится к большому числу тем (например, энциклопедия, журнал, сборник статей), то в задачах тематического поиска или классификации документов его имеет смысл разбивать на части, более однородные по тематике.

Если термин относится к большому числу тем, то, скорее всего, это общеупотребительное слово (стоп-слово), бесполезное для определения тематики.

Частотные (выборочные) оценки вероятностей. Вероятности, связанные с наблюдаемыми переменными d и w , можно оценивать по выборке как частоты (здесь и далее выборочные оценки вероятностей p будем обозначать через \hat{p}):

$$\hat{p}(d, w) = \frac{n_{dw}}{n}, \quad \hat{p}(d) = \frac{n_d}{n}, \quad \hat{p}(w) = \frac{n_w}{n}, \quad \hat{p}(w | d) = \frac{n_{dw}}{n_d}, \quad (1.3)$$

n_{dw} — число вхождений термина w в документ d ;

$n_d = \sum_{w \in W} n_{dw}$ — длина документа d в терминах;

$n_w = \sum_{d \in D} n_{dw}$ — число вхождений термина w во все документы коллекции;

$n = \sum_{d \in D} \sum_{w \in W} n_{dw}$ — длина коллекции в терминах.

Вероятности, связанные со скрытой переменной t , также можно оценивать как частоты, если рассматривать коллекцию документов как выборку троек (d, w, t) :

$$\hat{p}(t) = \frac{n_t}{n}, \quad \hat{p}(w | t) = \frac{n_{wt}}{n_t}, \quad \hat{p}(t | d) = \frac{n_{dt}}{n_d}, \quad \hat{p}(t | d, w) = \frac{n_{dwt}}{n_{dw}}, \quad (1.4)$$

n_{dwt} — число троек, в которых термин w документа d связан с темой t ;

$n_{dt} = \sum_{w \in W} n_{dwt}$ — число троек, в которых термин документа d связан с темой t ;

$n_{wt} = \sum_{d \in D} n_{dwt}$ — число троек, в которых термин w связан с темой t ;

$n_t = \sum_{d \in D} \sum_{w \in W} n_{dwt}$ — число троек, связанных с темой t .

В пределе $n \rightarrow \infty$ частотные оценки $\hat{p}(\cdot)$, определяемые формулами (1.3)–(1.4), стремятся к соответствующим вероятностям $p(\cdot)$, согласно закону больших чисел. Частотная интерпретация даёт ясное понимание всех условных вероятностей, которые будут использоваться в дальнейшем.

Связь с задачами неотрицательного матричного разложения. Если число тем $|T|$ много меньше числа документов $|D|$ и числа терминов $|W|$, то равенство (1.2) можно понимать как задачу приближённого представления заданной матрицы частот $F = (\hat{p}_{wd})_{W \times D}$, где $\hat{p}_{wd} = \hat{p}(w | d) = n_{dw}/n_d$, в виде произведения $F \approx \Phi\Theta$ двух неизвестных матриц меньшего размера — *матрицы тем* $\Phi = (\varphi_{wt})_{W \times T}$, $\varphi_{wt} = p(w | t)$ и *матрицы документов* $\Theta = (\theta_{td})_{T \times D}$, $\theta_{td} = p(t | d)$.

Одно из таких представлений строится из $|T|$ главных компонент сингулярного разложения матрицы F и является решением задачи наименьших квадратов

$$\sum_{d \in D} \sum_{w \in W} (\hat{p}_{wd} - p(w | d))^2 = \sum_{d \in D} \sum_{w \in W} \left(\hat{p}_{wd} - \sum_{t \in T} \varphi_{wt} \theta_{td} \right)^2 = \|F - \Phi\Theta\|^2 \rightarrow \min_{\Theta, \Phi}. \quad (1.5)$$

Сингулярное разложение имеет массу приложений в анализе данных, но для тематического моделирования оно плохо подходит. Во-первых, столбцы получаемых матриц Θ и Φ не удовлетворяют условиям неотрицательности и нормировки, поэтому их нельзя интерпретировать как распределения. Во-вторых, квадратичная функция потерь не чувствительна к малым различиям «хвостов» распределений, из-за которых их статистические свойства могут различаться существенно.

В вероятностном тематическом моделировании вместо принципа наименьших квадратов используется принцип максимума правдоподобия. Он также приводит к задаче матричного разложения вида (1.5), только вместо евклидовой нормы используется взвешенная дивергенция Кульбака–Лейблера.

§1.2 Предварительная обработка данных

Понятие «термина» может изменяться в зависимости от целей построения тематической модели и таких особенностей задачи, как язык документов, средняя длина документов, тематика коллекции.

Лемматизация и стемминг. При построении тематической модели нет смысла различать формы (склонения, спряжения) одного и того же слова. Это приведёт к неоправданному разрастанию словаря, дроблению статистики, увеличению ресурсоёмкости и снижению качества модели.

Лемматизация — это приведение каждого слова в документе к его нормальной форме. В русском языке нормальными формами считаются: для существительных — именительный падеж, единственное число; для прилагательных — именительный падеж, единственное число, мужской род; для глаголов, причастий, деепричастий — глагол в инфинитиве. Существуют специальные программы — *лемматизаторы* (lemmatizer), обычно основанные на явном хранении грамматического словаря со всеми формами слов. Недостатком лемматизации является трудоёмкость составления словарей, и, как следствие, их неполнота, особенно по части специальной терминологии и неологизмов, которые во многих приложениях как раз и представляют наибольший интерес.

[ссылка на рекомендуемые русский и английский лемматизаторы](#)

ToDo¹

Стемминг — это более простая технология, которая состоит в отбрасывании изменяемых частей слов, главным образом, окончаний. Она не требует хранения словаря всех слов и основана на правилах морфологии языка. Недостатком стемминга является большее число ошибок. Стемминг хорошо подходит для английского языка, но хуже подходит для русского.

[ссылка на рекомендуемые русский и английский стеммеры](#)

ToDo²

Отбрасывание стоп-слов. Слова, встречающиеся во многих текстах различной тематики, бесполезны для тематического моделирования, и могут быть отброшены. К ним относятся предлоги, союзы, числительные, местоимения, некоторые глаголы, прилагательные и наречия. Число таких слов обычно варьируется в пределах нескольких сотен. Их отбрасывание почти не влияет на длину словаря, но может приводить к заметному сокращению длины некоторых текстов.

Отбрасывание редких слов. Слова, встречающиеся в длинном документе слишком редко, например, только один раз, также можно отбрасывать, полагая, что данное слово не характеризует тематику данного документа. При обработке коллекций коротких новостных сообщений этот приём лучше не использовать.

Выделение ключевых фраз. При обработке коллекций научных, юридических или других специальных текстов вместо отдельных слов выделяют *ключевые фразы* — словосочетания, являющиеся устойчивыми оборотами или терминами в данной предметной области. Это отдельная довольно сложная задача, для решения которой используются тезаурусы, составленные экспертами [2], либо методы машинного обучения [26, 45], при этом для формирования обучающих выборок всё равно приходится привлекать экспертов.

Далее будем полагать, что словарь W получен в результате предварительной обработки всех документов коллекции D и может содержать как отдельные слова, так и ключевые фразы. Элементы словаря $w \in W$ будем называть «терминами».

§1.3 Метод максимума правдоподобия

Для оценивания параметров Θ, Φ тематической модели по коллекции документов D будем максимизировать правдоподобие (плотность распределения) выборки:

$$p(D; \Theta, \Phi) = C \prod_{d \in D} \prod_{w \in d} p(d, w)^{n_{dw}} = \prod_{d \in D} \prod_{w \in d} p(w | d)^{n_{dw}} \underbrace{C p(d)^{n_{dw}}}_{\text{const}} \rightarrow \max_{\Theta, \Phi}$$

где C — нормировочный множитель мультиномиального распределения, зависящий только от чисел n_{dw} . Отбросим множители C и $p(d)$, не влияющие на положение точки максимума, подставим выражение для $p(w | d)$ из (1.2) и воспользуемся обозначениями $\theta_{td} = p(t | d)$, $\varphi_{wt} = p(w | t)$. Прологарифмировав правдоподобие, получим задачу максимизации

$$L(D; \Theta, \Phi) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \varphi_{wt} \theta_{td} \rightarrow \max_{\Theta, \Phi} \quad (1.6)$$

при ограничениях неотрицательности $\theta_{td} \geq 0$, $\varphi_{wt} \geq 0$ и нормировки

$$\sum_{w \in W} \varphi_{wt} = 1, \quad \sum_{t \in T} \theta_{td} = 1.$$

Заметим, что максимизация (1.6) эквивалентна минимизации взвешенной суммы расстояний Кульбака–Лейблера $\text{KL}(\hat{p}||p) = \sum_w \hat{p}_{wd} \ln \frac{\hat{p}_{wd}}{p(w|d)}$ между эмпирическими распределениями \hat{p}_{wd} и модельными $p(w|d)$ по всем документам d из D :

$$\sum_{d \in D} n_d \text{KL} \left(\frac{n_{dw}}{n_d} \parallel \sum_{t \in T} \varphi_{wt} \theta_{td} \right) \rightarrow \min_{\Theta, \Phi},$$

где весом документа d является его длина n_d . Если веса n_d убрать, то все документы будут искусственно приведены к одинаковой длине. Такая модификация функционала качества бывает полезна при моделировании коллекций, содержащих документы одинаковой важности, но существенно разной длины.

§1.4 Внутренние оценки качества тематических моделей

Оценивание качества тематических моделей является нетривиальной проблемой. В отличие от задач классификации или регрессии здесь нет чёткого понятия «ошибки» или «потери». Стандартные критерии качества кластеризации типа средних внутрикластерных или межкластерных расстояний или их отношений плохо подходят для оценивания «мягкой» совместной кластеризации документов и терминов.

Наиболее распространённым критерием является *перплексия* (perplexity), используемая для оценивания моделей языка в компьютерной лингвистике. Это мера несоответствия или «удивлённости» модели $p(w|d)$ терминам w , наблюдаемым в документах d коллекции D , определяемая через логарифм правдоподобия (1.6):

$$\mathcal{P}(D; p) = \exp \left(-\frac{1}{n} L(D; \Theta, \Phi) \right) = \exp \left(-\frac{1}{n} \sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w|d) \right). \quad (1.7)$$

Чем меньше эта величина, тем лучше модель p предсказывает появление терминов w в документах d коллекции D .

Интерпретация перплексии. Если термины w порождаются из равномерного распределения $p(w) = 1/V$ на словаре мощности V , то перплексия модели p на таком тексте сходится к V с ростом его длины. Чем сильнее распределение p отличается от равномерного, тем меньше перплексия. Чем сильнее модель p отличается от генерирующего распределения, тем больше перплексия. В нашем случае в (1.7) используются условные вероятности терминов $p(w|d)$, и интерпретация немного другая: если каждый документ генерируется из V равновероятных терминов (возможно, различных в разных документах), то перплексия сходится к V . Опять-таки, чем сильнее распределение отличается от равномерного, тем меньше перплексия.

Чтобы сравнение перплексии двух коллекций было корректным, необходимо, чтобы они имели один и тот же словарь. ToDo³

Чтобы перплексия была характеристикой только качества модели, необходимо вводить нормировки, чтобы длины документов и эффективная мощность словаря не влияли на перплексию. ToDo⁴

Перплексия контрольной выборки. Обозначим через $p_D(w | d)$ модель, построенную по обучающей коллекции документов D . Перплексия обучающей выборки $\mathcal{P}(D; p_D)$ является оптимистично смещённой (заниженной) характеристикой качества модели из-за эффекта переобучения. Обобщающую способность модели принято оценивать *перплексией контрольной выборки* (hold-out perplexity) $\mathcal{P}(D'; p_D)$.

Вопрос о том, как разделить исходную коллекцию на обучение D и контроль D' , не тривиален. К сожалению, детали этой процедуры во многих статьях опускаются. В [6] предлагается разделять все документы на обучающие и контрольные случайным образом в пропорции 9 : 1. Однако в силу гипотез «мешка слов» и «мешка документов» более корректным было бы случайное разбиение каждого документа на обучающую и контрольную части. С другой стороны, во многих приложениях важно проверить способность тематической модели хорошо описывать новые документы.

Новые документы порождают две проблемы: во-первых, для них необходимо оценивать θ_{td} ; во-вторых, они могут содержать новые термины w , для которых придётся оценивать также φ_{wt} , увеличивать размерность векторов $\varphi_t = (\varphi_{wt})_{w \in W}$ и перенормировать их. Такая процедура оценивания модели частично включает в себя процедуру обучения, в результате чего оценка качества снова может оказаться оптимистично смещённой.

Частичное решение этой проблемы предлагается в [5]. После обучения модели p_D векторы φ_t фиксируются, векторы θ_d контрольных документов $d \in D'$ оцениваются по первой половине каждого документа, по вторым половинам вычисляется контрольная перплексия. Что такое «половина», не уточняется. Простое разрезание текста на две части может приводить к смещённым оценкам. Например, научные статьи обычно начинаются с введения и обзора, использующих общую терминологию, затем идёт изложение частных результатов. Если в коллекции много таких текстов, то оценка окажется пессимистично смещённой. Противоположный пример неслучайного разбиения текста — когда число вхождений каждого термина n_{dw} делится ровно пополам между обучающей и контрольной выборками. В таком случае обучающая и контрольная половины документа будут неразличимы для тематической модели, и оценка окажется оптимистично смещённой.

В наших экспериментах последовательность терминов $\{w_1, \dots, w_{n_d}\}$ каждого контрольного документа $d \in D'$ после случайной перестановки разбивается на две части равной длины. Новые слова, попадающие во вторую часть, игнорируются.

Ещё один выход — робастные модели. Новые редкие слова считаются шумом, описываются униграммной моделью и почти не дают вклада в контрольную перплексию. Робастную модель трудно удивить новыми словами, т.к. она трактует их как шум. ToDo⁵

Более сложные процедуры несмещённого оценивания правдоподобия предложены в [37] и улучшены в [7]. Они имеют трудоёмкость, квадратичную по длине документа, и в процессе оценивания используют ту же тематическую модель, качество которой оценивается. Эти недостатки несколько ограничивают их применимость.

Эксперименты на модельных данных. Алгоритм 1.1 можно использовать для генерации модельных данных по заданным распределениям $p(w | t)$ и $p(t | d)$. Это крайне полезно на стадии тестирования методов обучения тематических моделей, решающих задачу (1.6). Хороший метод должен быть способен восстановить по данным ту самую модель, которая эти данные породила. Модельные данные можно

генерировать различной длины n ; можно добавлять в них шум — случайные пары (d_i, w_i) из распределения, заведомо плохо приближаемого моделью (1.2); можно задавать распределения $p(w | t)$, $p(t | d)$ более различными или более похожими, тем самым делая задачу восстановления модели более лёгкой или более трудной; задавать различное число тем $|T|$, а восстанавливать модель при другом числе тем, либо пытаться его определить. Эксперименты с варьированием модели данных позволяют исследовать устойчивость метода и узнать границы его применимости. Только в случае модельных данных известно, какая тема t_i на самом деле связана с каждой парой (d_i, w_i) , что позволяет оценивать качество восстановления модели по данным как долю правильно угаданных тем или как расстояние между восстановленными и истинными распределениями $p(w | t)$, $p(t | d)$.

Показать эксперименты на модельных данных

ToDo⁶

2 Вероятностный латентный семантический анализ

Вероятностный латентный семантический анализ (probabilistic latent semantic analysis, PLSA) был предложен Томасом Хоффманом в [17].

Вероятностная модель появления пары «документ–термин» (d, w) записывается тремя эквивалентными способами:

$$p(d, w) = \sum_{t \in T} p(t)p(w | t)p(d | t) = \sum_{t \in T} p(d)p(w | t)p(t | d) = \sum_{t \in T} p(w)p(t | w)p(d | t),$$

где $p(t)$ — распределение тем во всей коллекции. Первое представление называется симметричным, второе и третье — несимметричными. Они приводят к немного разным итерационным процессам обучения тематической модели. Сейчас возьмём за основу второе представление, совпадающее с (1.2).

§2.1 EM-алгоритм

Для решения задачи (1.6) в PLSA применяется итерационный процесс, в котором каждая итерация состоит из двух шагов — E (expectation) и M (maximization) [12]. Перед первой итерацией выбирается начальное приближение параметров φ_{wt} , θ_{td} .

На E-шаге по текущим значениям параметров φ_{wt} , θ_{td} с помощью формулы Байеса вычисляются условные вероятности $p(t | d, w)$ всех тем $t \in T$ для каждого термина $w \in d$ в каждом документе d :

$$H_{dwt} = p(t | d, w) = \frac{p(w | t)p(t | d)}{p(w | d)} = \frac{\varphi_{wt}\theta_{td}}{\sum_{s \in T} \varphi_{ws}\theta_{sd}}. \quad (2.1)$$

На M-шаге, наоборот, по условным вероятностям тем H_{dwt} вычисляется новое приближение параметров φ_{wt} , θ_{td} . Это легко сделать, если заметить, что величина $\hat{n}_{dwt} = n_{dw}H_{dwt}$ оценивает (не обязательно целое) число n_{dwt} вхождений термина w в документ d , связанных с темой t . Просуммировав \hat{n}_{dwt} по документам d и по терминам w , получим оценки \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t , и через них, согласно (1.4), — частотные оценки

условных вероятностей φ_{wt} , θ_{td} :

$$\varphi_{wt} = \frac{\hat{n}_{wt}}{\hat{n}_t}, \quad \hat{n}_t = \sum_{w \in W} \hat{n}_{wt}, \quad \hat{n}_{wt} = \sum_{d \in D} n_{dw} H_{dwt}. \quad (2.2)$$

$$\theta_{td} = \frac{\hat{n}_{dt}}{\hat{n}_d}, \quad \hat{n}_d = \sum_{t \in T} \hat{n}_{dt}, \quad \hat{n}_{dt} = \sum_{w \in d} n_{dw} H_{dwt}. \quad (2.3)$$

Эти простые, но не вполне строгие рассуждения поясняют суть EM-алгоритма. Покажем теперь, что оценки (2.2)–(2.3) действительно являются решением задачи максимизации правдоподобия (1.6) при фиксированных H_{dwt} . Запишем лагранжиан задачи (1.6) при ограничениях нормировки, проигнорировав ограничения неотрицательности (позже убедимся, что решение действительно неотрицательно):

$$\mathcal{L}(\Theta, \Phi) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \underbrace{\sum_{t \in T} \varphi_{wt} \theta_{td}}_{p(w|d)} - \sum_{t \in T} \lambda_t \left(\sum_{w \in W} \varphi_{wt} - 1 \right) - \sum_{d \in D} \mu_d \left(\sum_{t \in T} \theta_{td} - 1 \right).$$

Продифференцировав лагранжиан по φ_{wt} и приравняв нулю производную, получим

$$\lambda_t = \sum_{d \in D} n_{dw} \frac{\theta_{td}}{p(w|d)}. \quad (2.4)$$

Домножим обе части этого равенства на φ_{wt} , просуммируем по всем терминам $w \in W$, применим условие нормировки вероятностей φ_{wt} в левой части и выделим переменную H_{dwt} в правой части. Получим

$$\lambda_t = \sum_{d \in D} \sum_{w \in W} n_{dw} H_{dwt}.$$

Снова домножим обе части (2.4) на φ_{wt} , выделим переменную H_{dwt} в правой части и выразим φ_{wt} из левой части, подставив уже известное выражение для λ_t . Получим

$$\varphi_{wt} = \frac{\sum_{d \in D} n_{dw} H_{dwt}}{\sum_{w' \in W} \sum_{d \in D} n_{dw'} H_{dw't}}.$$

Обозначив числитель через \hat{n}_{wt} , получим (2.2). Прделав аналогичные действия с производной лагранжиана по θ_{td} , получим (2.3).

Эффективность EM-алгоритма по времени и по памяти. Число операций растёт линейно по длине коллекции n , числу тем T и числу итераций.

Перебор всех терминов w во всех документах d можно организовать очень эффективно, если хранить каждый документ d в виде последовательности пар (w, n_{dw}) .

Рациональный EM-алгоритм. Вычисление переменных \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t на M-шаге требует однократного прохода всей коллекции в цикле по всем документам $d \in D$ и всем терминам $w \in d$. Внутри этого цикла переменные H_{dwt} можно вычислять непосредственно в тот момент, когда они понадобятся. От этого результат алгоритма не изменяется, E-шаг встраивается внутрь M-шага без дополнительных вычислительных затрат, отпадает необходимость хранения трёхмерной матрицы H_{dwt} . Заметим также, что переменную \hat{n}_d можно не вычислять, поскольку $\hat{n}_d = n_d$. Этот вариант реализации EM-алгоритма будем называть *рациональным*; он показан в Алгоритме 2.1.

Алгоритм 2.1. PLSA-EM: рациональный EM-алгоритм для модели PLSA.

Вход: коллекция документов D , число тем $|T|$, начальные приближения Θ , Φ ;

Выход: распределения Θ и Φ ;

```

1 повторять
2   обнулить  $\hat{n}_{wt}$ ,  $\hat{n}_{dt}$ ,  $\hat{n}_t$  для всех  $d \in D$ ,  $w \in W$ ,  $t \in T$ ;
3   для всех  $d \in D$ ,  $w \in d$ 
4      $Z := \sum_{t \in T} \varphi_{wt} \theta_{td}$ ;
5     для всех  $t \in T$  таких, что  $\varphi_{wt} \theta_{td} > 0$ 
6      $\lceil$  увеличить  $\hat{n}_{wt}$ ,  $\hat{n}_{dt}$ ,  $\hat{n}_t$  на  $\delta = n_{dw} \varphi_{wt} \theta_{td} / Z$ ;
7    $\varphi_{wt} := \hat{n}_{wt} / \hat{n}_t$  для всех  $w \in W$ ,  $t \in T$ ;
8    $\theta_{td} := \hat{n}_{dt} / n_d$  для всех  $d \in D$ ,  $t \in T$ ;
9 пока  $\Theta$  и  $\Phi$  не сойдутся;
  
```

Проблема разреженности. Если начальные приближения θ_{td} и φ_{wt} положительны, то и после каждой итерации они будут оставаться положительными, несмотря на то, что ограничение неотрицательности было проигнорировано в ходе решения. И, наоборот, если $\theta_{td} = 0$ (тема t не представлена в документе d) или если $\varphi_{wt} = 0$ (термин w не относится к теме t), то нулевое значение будет сохраняться на протяжении всех итераций. Таким образом, в PLSA структура разреженности распределений не оптимизируется, а задаётся через начальное приближение. В то же время, использование разреженных матриц для хранения переменных \hat{n}_{wt} , \hat{n}_{dt} , θ_{td} , φ_{wt} могло бы дать существенную экономию памяти.

Эксперимент: принудительное разреживание портит модель, если его делать с первых итераций. Надо дождаться сходимости, когда правильно определится подмножество малых вероятностей. ToDo⁷

§2.2 Обобщённый EM-алгоритм

В EM-алгоритме нет необходимости сверхточно решать задачу максимизации правдоподобия на M-шаге. Достаточно ещё немного приблизиться к точке максимума правдоподобия и снова выполнить E-шаг. Это связано с тем, что сам функционал правдоподобия известен не точно — он зависит от приближённых значений H_{dwt} , полученных на E-шаге. EM-алгоритм с сокращённым M-шагом называется *обобщённым EM-алгоритмом* (generalized EM-algorithm, GEM). Для него справедливы те же доказательства сходимости, что и для основного варианта EM-алгоритма [12].

В случае PLSA сокращение M-шага сводится к более частому обновлению параметров θ_{td} и φ_{wt} по значениям счётчиков \hat{n}_{wt} и \hat{n}_{dt} . В Алгоритме 2.1 это происходит после каждого просмотра всей коллекции. Обновления можно делать после обработки каждого документа или после заданного числа обработанных пар (d, w) или даже после каждой пары. На больших коллекциях частые обновления повышают скорость сходимости. В Алгоритме 2.2 выбор условия обновления на шаге 9 оставлен на усмотрение разработчика.

Алгоритм 2.2. PLSA-GEM: обобщённый EM-алгоритм для модели PLSA.

Вход: коллекция документов D , число тем $|T|$, начальные приближения Θ , Φ ;

Выход: распределения Θ и Φ ;

1 обнулить \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t , \hat{n}_d , n_{dwt} для всех $d \in D$, $w \in W$, $t \in T$;

2 **повторять**

3 **для всех** $d \in D$, $w \in d$

4 $Z := \sum_{t \in T} \varphi_{wt} \theta_{td}$;

5 **для всех** $t \in T$ таких, что $n_{dwt} > 0$ или $\varphi_{wt} \theta_{td} > 0$

6 $\delta := n_{dw} \varphi_{wt} \theta_{td} / Z$;

7 увеличить \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t , \hat{n}_d на $(\delta - n_{dwt})$;

8 $n_{dwt} := \delta$;

9 **если** пора обновить параметры Φ , Θ **то**

10 $\varphi_{wt} := \hat{n}_{wt} / \hat{n}_t$ для всех $w \in W$, $t \in T$ таких, что \hat{n}_{wt} изменился;

11 $\theta_{td} := \hat{n}_{dt} / \hat{n}_d$ для всех $d \in D$, $t \in T$ таких, что \hat{n}_{dt} изменился;

12 **пока** Θ и Φ не сойдутся;

На первой итерации (т. е. при первом проходе коллекции) частые обновления не делаются, чтобы в счётчиках накопилась информация по всей коллекции. В противном случае оценки параметров θ_{td} и φ_{wt} по начальному фрагменту выборки могут оказаться хуже начального приближения. Начиная со второй итерации, для каждой пары (d, w) из счётчиков \hat{n}_{wt} и \hat{n}_{dt} вычитается n_{dwt} — то самое значение δ , которое было к ним прибавлено при обработке пары (d, w) на предыдущей итерации. Таким образом, счётчики \hat{n}_{wt} и \hat{n}_{dt} всегда содержат результат последнего однократного прохода всей коллекции.

Необходимость хранения трёхмерной матрицы n_{dwt} делает Алгоритм 2.2 неприменимым к большим коллекциям. Этот недостаток устраняется путём реорганизации итераций, либо применением сэмплирования. Далее рассматриваются оба способа.

Эксперимент: частота обновления влияет на эффективность но не влияет на качество. Лучше всего обновлять после каждого слова. При этом можно вообще отказаться от хранения матриц тета и фи.

ToDo⁸

§2.3 Онлайновый EM-алгоритм

На больших коллекциях Алгоритмы 2.1 и 2.2 могут сходиться очень медленно. Причина в том, что за однократный проход по всем документам коллекции оценки распределений терминов в темах $\varphi_{wt} = \hat{n}_{wt} / \hat{n}_t$ уточняются огромное число раз и успевают сойтись, в то же время распределения тем в документах θ_d проходят лишь одну итерацию. На начальных итерациях, пока распределения θ_d не сошлись, вычислительный ресурс тратится впустую на достижение сходимости φ_t к приближениям, далёким от оптимальных. Суть этой проблемы в том, что параметры θ_{td} привязаны к отдельным документам d , а параметры φ_{wt} — ко всей коллекции.

Проблема решается реорганизацией шагов итерационного процесса. Проход каждого документа $d \in D$ производится несколько раз подряд. На каждом проходе

Алгоритм 2.3. PLSA-BatchEM: пакетный EM-алгоритм для модели PLSA.

Вход: коллекция документов D , число тем $|T|$;

Выход: распределения Θ и Φ ;

```

1 инициализировать  $\varphi_{wt}$  для всех  $w \in W, t \in T$ ;
2 повторять
3    $\hat{n}_{wt} := 0; \hat{n}_t := 0$  для всех  $w \in W, t \in T$ ;
4   для всех  $d \in D$ 
5     инициализировать  $\theta_{td}$  для всех  $t \in T$ ;
6     повторять
7        $Z_w := \sum_{t \in T} \varphi_{wt} \theta_{td}$  для всех  $w \in d$ ;
8        $\theta_{td} := \frac{1}{n_d} \sum_{w \in d} n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $t \in T$ ;
9     пока  $\theta_d$  не сойдётся;
10    увеличить  $\hat{n}_{wt}, \hat{n}_t$  на  $n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $w \in d, t \in T$ ;
11    $\varphi_{wt} := \hat{n}_{wt} / \hat{n}_t$  для всех  $w \in W, t \in T$ ;
12 пока  $\Phi$  не сойдётся;
```

документа выполняется E-шаг и обновляется распределение θ_d . Обновление распределений φ_t производится после каждого прохода коллекции. В результате распределения φ_t и θ_d сходятся более согласованно.

Реорганизация позволяет отказаться от трёхмерных массивов. В псевдокоде Алгоритма 2.3 используются переменные $H_{wt} = p(t | d, w)$ вместо H_{dwt} , чтобы подчеркнуть, что по окончании обработки документа d эти данные уже не нужны, и двумерный массив $(H_{wt})_{W \times T}$ можно использовать для обработки следующего документа.

Можно также отказаться от двумерных массивов, размер которых пропорционален $|D|$, что позволит обрабатывать очень большие коллекции документов. Вероятности θ_{td} всех тем $t \in T$ документа d не нужны по окончании обработки документа d , поэтому двумерный массив $(\theta_{td})_{T \times D}$ можно заменить одномерным $(\theta_t)_T$.

Хранение двумерного массива $(\theta_{td})_{T \times D}$ всё же имеет смысл, если размер коллекции $|D|$ относительно невелик, и на шаге 5 инициализация θ_{td} производится только во время первого прохода коллекции, а при последующих проходах используется текущая оценка θ_{td} , оставшаяся с предыдущего прохода. Хорошее начальное приближение обеспечивает сходимость θ_d за меньшее число итераций.

Скорость сходимости зависит от тщательности подбора числа итераций на внутреннем цикле по документу и внешнем цикле по коллекции. На начальных итерациях внешнего цикла можно делать меньше итераций внутреннего цикла, поскольку нет смысла добиваться сходимости θ_d , пока распределения φ_t далеки от оптимальных.

Ещё одна идея ускорения сходимости состоит в том, чтобы начальные итерации провести не по всей коллекции, а по случайному подмножеству (пакету) документов $D' \subseteq D$. Если коллекция имеет избыточный размер, то для получения хорошего приближения Φ достаточно будет просмотреть относительно небольшую её часть.

Алгоритм 2.3 назван *пакетным* (batch algorithm), так как он может обрабатывать коллекцию по частям. Развитие этой идеи приводит к онлайн-алгоритму

Алгоритм 2.4. PLSA-OEM: онлайнный EM-алгоритм для модели PLSA.

Вход: коллекция документов D , число тем $|T|$;

Выход: распределения Θ и Φ ;

```

1 инициализировать  $\varphi_{wt}$  для всех  $w \in W, t \in T$ ;
2  $\hat{n}_{wt} := 0, \hat{n}_t := 0$  для всех  $w \in W, t \in T$ ;
3 для всех пакетов  $D_j, j = 1, \dots, J$ 
4   повторять
5      $\tilde{n}_{wt} := 0, \tilde{n}_t := 0$  для всех  $w \in W, t \in T$ ;
6     для всех  $d \in D_j$ 
7       инициализировать  $\theta_{td}$  для всех  $t \in T$ ;
8       повторять
9          $Z_w := \sum_{t \in T} \varphi_{wt} \theta_{td}$  для всех  $w \in d$ ;
10         $\theta_{td} := \frac{1}{n_d} \sum_{w \in d} n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $t \in T$ ;
11        пока  $\theta_d$  не сойдётся;
12        увеличить  $\tilde{n}_{wt}, \tilde{n}_t$  на  $n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $w \in d, t \in T$ ;
13         $\varphi_{wt} := \frac{\rho_j \hat{n}_{wt} + \tilde{n}_{wt}}{\rho_j \hat{n}_t + \tilde{n}_t}$  для всех  $w \in W, t \in T$  таких, что  $\tilde{n}_{wt} > 0$ ;
14        пока  $\Phi$  не сойдётся;
15         $\hat{n}_{wt} := \rho_j \hat{n}_{wt} + \tilde{n}_{wt}$  для всех  $w \in W, t \in T$ ;
16         $\hat{n}_t := \rho_j \hat{n}_t + \tilde{n}_t$  для всех  $t \in T$ ;

```

му [16], одному из самых быстрых в тематическом моделировании. Он реализован в библиотеке онлайнных алгоритмов Vowpal Wabbit Джона Лэнгфорда.

Онлайнный алгоритм. В машинном обучении *онлайнными* принято называть алгоритмы, способные адекватно настраивать параметры модели за один проход по выборке. Онлайнные алгоритмы используются для обработки потоковых данных. Во многих приложениях тематического моделирования коллекция документов пополняется динамически, и требуется обновлять модель, не обрабатывая заново всю коллекцию. Обновление модели предполагает вычисление распределения $\theta_{td} = p(t | d)$ для нового документа d и уточнение распределений $\varphi_{wt} = p(w | t)$ для всех тем $t \in T$, имеющих ненулевые вероятности для слов документа d . Если коллекция уже имеет большой объём, то добавление документа не должно сильно повлиять на Φ . Чем больше коллекция, тем лучше текущее приближение Φ , и тем меньше итераций потребуется для добавления нового документа.

Стратегии ускорения сходимости. Первый пакет — особенный! Для первого пакета: 1) число итераций θ_d ограничить сверху числом прошедших итераций Φ плюс несколько.

2) использовать двумерный массив θ_{td} и инициализировать его только при первом проходе.

Онлайнный Алгоритм 2.4 является модификацией пакетного Алгоритма 2.3. Теперь вся коллекция разбивается на пакеты документов $D_1, D_2, \dots, D_j, \dots$. Способ разбиения остаётся на усмотрение разработчика, в частности, пакеты могут пере-

ToDo⁹

секаться либо не пересекаться, просматриваться по одному разу либо многократно, выбираться случайно, по времени поступления, либо по времени публикации документа, и т. д. Размер первого пакета $|D_1|$ должен быть достаточным для получения распределений φ_{wt} с приемлемой точностью. Обработка каждого пакета производится пакетным Алгоритмом 2.3 при фиксированных φ_{wt} . Затем счётчики \tilde{n}_{wt} , вычисленные по обработанному пакету документов, складываются со счётчиками \hat{n}_{wt} , вычисленными по всем предыдущим пакетам.

Если делается много проходов по коллекции или если значимость пакетов убывает по мере поступления новых, то старые счётчики домножаются на параметр $\rho_j \in (0, 1]$, управляющий скоростью забывания старых оценок:

$$\hat{n}_{wt} := \rho_j \hat{n}_{wt} + \tilde{n}_{wt}.$$

Фактически, при поступлении каждого нового пакета документов D_j частоты слов во всех старых документах уменьшаются в ρ_j^{-1} раз. При $\rho_j = 1$ забывания нет и φ_{wt} вычисляются как обычные частотные оценки условных вероятностей.

§2.4 Стохастический EM-алгоритм

В Алгоритме 2.2 для каждой пары (d, w) происходит распределение n_{dw} входящий термина w в документ d между всеми $|T|$ темами пропорционально вероятностям $p(t | d, w)$. При этом приходится хранить массив значений n_{dwt} для всех тем $t \in T$. Расход памяти объёма $O(n|T|)$ может оказаться неприемлемым даже при небольшом числе тем. В то же время, согласно гипотезе разреженности, употребление термина w в документе d связано, скорее всего, с небольшим числом тем.

Можно было бы оставлять только несколько наибольших значений n_{dwt} на каждом шаге. Однако эксперименты показывают, что эта эвристика приводит к накоплению систематической ошибки и смещению модели.

Эксперимент по тах-разреживанию для PLSA.

ToDo¹⁰

Проблема разреживания условного распределения $p(t | d, w)$ адекватно решается с помощью стохастического EM-алгоритма (stochastic EM-algorithm, SEM) [8]. Распределение скрытой переменной t , вычисленное на E-шаге, не используется непосредственно на M-шаге. Вместо этого из него сэмплируется искусственная выборка, по этой выборке вычисляется эмпирическое распределение, и оно уже используется в формулах M-шага. Это позволяет упростить задачу M-шага, сохранив свойства несмещённости оценок и сходимости EM-алгоритма. Размер сэмплируемой выборки является параметром метода.

В случае PLSA реализация SEM сводится к следующему: для каждой пары (d, w) сэмплируются s случайных тем t_{dwi} , $i = 1, \dots, s$ из распределения $p(t | d, w)$, возможно, повторяющихся. В формулах M-шага вместо распределения $p(t | d, w)$ используется его несмещённая эмпирическая оценка:

$$\hat{p}(t | d, w) = \frac{1}{s} \sum_{i=1}^s [t_{dwi} = t]. \quad (2.5)$$

Модификация Алгоритма 2.2, трансформирующая его в стохастический обобщённый EM-алгоритм (PLSA-SEM), состоит из трёх изменений:

- 1) перед шагом 5 сэмплируется s тем t_{dwi} , $i = 1, \dots, s$ из $p(t | d, w)$;
- 2) на шаге 5 цикл по всем $t \in T$ заменяется циклом по $t = t_{dwi}$, $i = 1, \dots, s$;
- 3) на шаге 6 вычисляется $\delta := n_{dw}/s$.

При $s = n_{dw}$ стохастический EM-алгоритм соответствует *сэмплированию Гиббса* [38], которое считается одним из основных методов обучения вероятностных тематических моделей. Однако эксперименты показывают, что параметр s можно брать намного меньше, от 1 до 5. Эта эвристика, названная *экономным сэмплированием* [1], приводит к разреживанию распределений $p(t | d, w)$ и существенной экономии вычислительного ресурса и памяти без потери качества тематической модели.

Эксперимент: зависимость контрольной перплексии от параметра разреживания s . ToDo¹¹

§2.5 Формирование начальных приближений

Начальные приближения φ_t и θ_d можно задавать нормированными случайными векторами из равномерного распределения.

Другая распространённая рекомендация — пройти по всей коллекции, выбрать для каждой пары (d, w) случайную тему t и вычислить частотные оценки (1.4) вероятностей φ_{wt} и θ_{td} для всех $d \in D$, $w \in W$, $t \in T$.

Инициализация с частичным обучением применяется в случаях, когда темы известны заранее и имеются дополнительные данные о привязке некоторых документов или терминов к темам. Учёт этих данных улучшает интерпретируемость тем.

Если известно, что документ d относится к подмножеству тем $T_d \subset T$, то в качестве начального θ_{td} можно взять равномерное распределение на этом подмножестве:

$$\theta_{td}^0 = \frac{1}{|T_d|} [t \in T_d]. \quad (2.6)$$

Если известно, что подмножество терминов $W_t \subset W$ относится к теме t , то в качестве начального φ_{wt} можно взять равномерное распределение на W_t :

$$\varphi_{wt}^0 = \frac{1}{|W_t|} [w \in W_t]. \quad (2.7)$$

Если известно, что подмножество документов $D_t \subset D$ относится к теме t , то можно взять эмпирическое распределение слов в объединённом документе:

$$\varphi_{wt}^0 = \frac{\sum_{d \in D_t} n_{dw}}{\sum_{d \in D_t} n_d}.$$

Если нет никакой априорной информации о связи документов с темами, то последнюю формулу можно применить к случайным подмножествам документов D_t . В [15] предлагается брать один случайный документ.

Инициализация Θ по Φ . Если для всех тем известны начальные приближения φ_{wt}^0 , то можно взять равномерное распределение $\theta_{td}^0 = 1/|T|$. Первая итерация EM-алгоритма даёт ещё одну интуитивно очевидную формулу инициализации:

$$\theta_{td} = \frac{1}{n_d} \sum_{w \in d} n_{dw} H_{dwt} = \sum_{w \in d} \frac{n_{dw}}{n_d} \frac{\varphi_{wt}}{\sum_s \varphi_{ws}}. \quad (2.8)$$

Здесь распределение тем в документе d оценивается путём усреднения распределений тем $p(t | w)$ по словам документа d , вычисленных по формуле Байеса.

Сглаживание. Если полученное начальное приближение φ_{wt}^0 или θ_{td}^0 содержит нулевые вероятности, то его необходимо сгладить, смешав с каким-нибудь неразрезанным распределением. Например, φ_{wt}^0 смешивается с эмпирическим распределением слов во всей коллекции и со случайным распределением $\rho(w)$, при некоторых значениях параметров смеси τ_1 и τ_2 :

$$\varphi_{wt} = (1 - \tau_1 - \tau_2)\varphi_{wt}^0 + \tau_1 n_w / n + \tau_2 \rho(w).$$

Эксперименты: сравнение нескольких способов задания начального приближения.

ToDo¹²

§2.6 Частичное обучение

В некоторых задачах классификации и категоризации текстов бывает известно, что какие-то документы или термины относятся или, наоборот, не относятся к некоторым темам. В таких случаях говорят о задачах с *частичным обучением* (semi-supervised learning). В EM-алгоритме нетрудно учесть данные такого типа. Рассмотрим для определённости Алгоритм 2.2. Модификации коснутся только правила пересчёта параметров θ_{td} и φ_{wt} на шагах 10–11. Те же модификации нетрудно сделать для пакетного, онлайн-ового или стохастического алгоритма.

Данные о релевантности. Данные о том, что документ d или термин w связан с темой t , можно использовать не только на этапе инициализации, но и далее на каждой итерации. Привязка к теме должна быть «мягкой», не исключающей связей с другими темами. Это делается с помощью стандартной модификация EM-алгоритма для задач с частичным обучением, которая используется в кластеризации и тематическом моделировании [25]. Она заключается в том, что на каждой итерации текущие приближения φ_{wt} и θ_{td} немного сдвигаются в сторону начального приближения:

$$\begin{aligned}\varphi_{wt} &:= \lambda \varphi_{wt}^0 + (1 - \lambda) \frac{n_{wt}}{n_t}; \\ \theta_{td} &:= \mu \theta_{td}^0 + (1 - \mu) \frac{n_{dt}}{n_d};\end{aligned}$$

где λ и μ — параметры, принимающие значения из отрезка $[0, 1]$; распределения φ_{wt}^0 и θ_{td}^0 вычисляются согласно (2.7) и (2.6).

Эксперимент А.Каца на модельных данных, подтверждающий, что достаточно привязать менее 10% документов, чтобы добиться однозначного соответствия исходной и восстановленной модели.

ToDo¹³

Продолжение эксперимента на реальных данных.

Гипотеза: учёт начальной информации на каждой итерации работает лучше, чем её учёт только на стадии инициализации. На модельных данных темы восстанавливаются точнее, возможно, требуя меньшего числа привязок. На реальных данных улучшается интерпретируемость тем.

Обязательно сделать подбор параметров λ и μ .

Виртуальные документы. Интерпретацию темы t фиксирует список ключевых терминов W_t . Он включается в коллекцию D как «виртуальный» документ d с заданными $n_{dw} = M_d[w \in W_t]$ и начальным распределением $\theta_{\tau d}^0 = [\tau = t]$, где $\tau \in T$. Увеличивая значение параметра M_d , можно повышать степень влияния виртуального документа d на тематическую модель.

Данные о нерелевантности. Чтобы учесть априорную информацию о том, что документ d или термин w не связан с темой t , достаточно обнулить параметр θ_{td} или φ_{wt} , обнулив соответствующие счётчики перед шагом 10:

если документ d не связан с темой t **то**

$$\begin{aligned}\hat{n}_d &:= \hat{n}_d - \hat{n}_{dt}; \\ \hat{n}_{dt} &:= 0;\end{aligned}$$

если термин w не связан с темой t **то**

$$\begin{aligned}\hat{n}_t &:= \hat{n}_t - \hat{n}_{wt}; \\ \hat{n}_{wt} &:= 0.\end{aligned}$$

Данные о переранжировании. Допустим, эксперты имеют возможность просмотреть список тем по любому документу d , ранжированный по убыванию частот \hat{n}_{dt} . Эксперт может перенести любую тему на то место в списке, которое он считает наиболее релевантным.

Эти данные, полученные от экспертов, легко учесть в EM-алгоритме. Чтобы тема t оказалась на k -м месте в списке тем документа d , достаточно сделать значение \hat{n}_{dt} немного большим k -го значения $\hat{n}_{dt}^{(k)}$ и скорректировать счётчик \hat{n}_d :

если тема t для документа d должна быть на k -м месте **то**

$$\begin{aligned}\hat{n}'_{dt} &:= \frac{1}{2}(\hat{n}_{dt}^{(k)} + \hat{n}_{dt}^{(k-1)})[k > 1] + \frac{1}{2}(3\hat{n}_{dt}^{(k)} - \hat{n}_{dt}^{(k+1)})[k = 1]; \\ \hat{n}_d &:= \hat{n}_d - \hat{n}_{dt} + \hat{n}'_{dt}; \\ \hat{n}_{dt} &:= \hat{n}'_{dt};\end{aligned}$$

Аналогичным образом возможно собрать и учесть данные о переранжировании терминов. Допустим, по теме t имеется список терминов, ранжированный по убыванию частот n_{wt} , и эксперт может проделать с ним аналогичную работу. Чтобы термин w оказался на k -м месте в списке терминов темы t , достаточно сделать значение \hat{n}_{wt} немного большим k -го значения $\hat{n}_{wt}^{(k)}$ и скорректировать счётчик \hat{n}_t :

если термин w для темы t должен быть на k -м месте **то**

$$\begin{aligned}\hat{n}'_{wt} &:= \frac{1}{2}(\hat{n}_{wt}^{(k)} + \hat{n}_{wt}^{(k-1)})[k > 1] + \frac{1}{2}(3\hat{n}_{wt}^{(k)} - \hat{n}_{wt}^{(k+1)})[k = 1]; \\ \hat{n}_t &:= \hat{n}_t - \hat{n}_{wt} + \hat{n}'_{wt}; \\ \hat{n}_{wt} &:= \hat{n}'_{wt};\end{aligned}$$

Существуют и другие способы задания априорной информации, для которых также можно адаптировать правила пересчёта распределений по счётчикам.

Эксперименты с виртуальными документами.

ToDo¹⁴

Стандартные методы SSL и сравнение с предложенным методом.

ToDo¹⁵

3 Латентное размещение Дирихле

Основным недостатком PLSA считается высокая размерность пространства параметров, вызывающая переобучение [6]. В задачах машинного обучения для сокращения размерности обычно используется либо *отбор признаков*, приводящий к уменьшению числа параметров, либо *регуляризация* — наложение дополнительных ограничений на параметры. В частности, *байесовская регуляризация* основана на введении априорного распределения вероятности в пространстве параметров.

Тематическая модель *латентного размещения Дирихле* (latent Dirichlet allocation, LDA) [6] основана на разложении (1.2) при дополнительном предположении, что векторы документов $\theta_d = (\theta_{td}) \in \mathbb{R}^{|T|}$ и векторы тем $\varphi_t = (\varphi_{wt}) \in \mathbb{R}^{|W|}$ порождаются распределениями Дирихле с параметрами $\alpha \in \mathbb{R}^{|T|}$ и $\beta \in \mathbb{R}^{|W|}$ соответственно:

$$\begin{aligned} \text{Dir}(\theta_d; \alpha) &= \frac{\Gamma(\alpha_0)}{\prod_t \Gamma(\alpha_t)} \prod_t \theta_{td}^{\alpha_t - 1}, \quad \alpha_t > 0, \quad \alpha_0 = \sum_t \alpha_t, \quad \theta_{td} > 0, \quad \sum_t \theta_{td} = 1; \\ \text{Dir}(\varphi_t; \beta) &= \frac{\Gamma(\beta_0)}{\prod_w \Gamma(\beta_w)} \prod_w \varphi_{wt}^{\beta_w - 1}, \quad \beta_w > 0, \quad \beta_0 = \sum_w \beta_w, \quad \varphi_{wt} > 0, \quad \sum_w \varphi_{wt} = 1. \end{aligned}$$

где $\Gamma(z)$ — гамма-функция. Считается, что распределение Дирихле хорошо подходит в качестве байесовского регуляризатора в задачах тематического моделирования.

Во-первых, это достаточно широкое параметрическое семейство распределений на единичном симплексе, то есть на множестве дискретных распределений. Если $\alpha_t = 1$ для всех t , то распределение Дирихле переходит в равномерное. Математическое ожидание и дисперсия t -й координаты вектора θ_d равны, соответственно,

$$\mathbf{E}\theta_{td} = \int \theta_{td} \text{Dir}(\theta_d; \alpha) d\theta_d = \frac{\alpha_t}{\alpha_0}, \quad \mathbf{D}\theta_{td} = \frac{\alpha_t(\alpha_0 - \alpha_t)}{\alpha_0^2(\alpha_0 + 1)}. \quad (3.1)$$

Векторный параметр α определяет степень разреженности векторов θ_d , порождаемых распределением $\text{Dir}(\theta; \alpha)$. Чем больше α_0 , тем сильнее векторы θ_d концентрируются вокруг вектора математического ожидания $\mathbf{E}\theta_d$. Чем меньше α_t , тем сильнее значения θ_{td} концентрируются вокруг нуля. Чем меньше α_0 , тем более разрежен вектор θ_d . Поэтому α_t называют *параметрами контраста*.

Во-вторых, модель LDA хорошо подходит для описания кластерных структур. Чем меньше α_0 и β_0 , тем сильнее разрежены распределения Дирихле, тем дальше отстоят друг от друга порождаемые ими векторы θ_d и φ_t , и тем чётче выражены тематические кластерные структуры в коллекции документов.

В-третьих, распределение Дирихле является сопряжённым к мультиномиальному, что упрощает вывод апостериорных оценок вероятностей θ_{td} и φ_{wt} .

Недостатком распределения Дирихле является отсутствие убедительных лингвистических обоснований. Второй недостаток в том, что параметры θ_{td} и φ_{wt} не могут обращаться в нуль, что противоречит гипотезе разреженности.

§3.1 Байесовский вывод

Рассмотрим процесс порождения документа d как выборки n_d пар тема–термин $X_d = \{(t_1, w_1), \dots, (t_{n_d}, w_{n_d})\}$. В каждой паре (t_i, w_i) тема t_i выбирается из дискретного распределения $p(t | d) = \theta_{td}$. Следовательно, вероятность встретить каждую

из тем t ровно n_{td} раз подчиняется мультиномиальному распределению:

$$p(X_d|\theta_d) = \frac{n_d!}{\prod_t n_{td}!} \prod_t \theta_{td}^{n_{td}}.$$

Распределение Дирихле является *сопряжённым* к мультиномиальному. Это означает, что при априорном распределении Дирихле $\theta_d \sim \text{Dir}(\theta; \alpha)$ апостериорное распределение вектора θ_d принадлежит тому же семейству распределений, но с другим значением параметра: $\theta_d|X_d \sim \text{Dir}(\theta; \alpha')$. Действительно, по формуле Байеса

$$p(\theta_d|X_d, \alpha) = \frac{p(X_d|\theta_d) \text{Dir}(\theta_d; \alpha)}{p(X_d)} = C \prod_t \theta_{td}^{n_{td}} \theta_{td}^{\alpha_t - 1} = \text{Dir}(\theta_d; \alpha'), \quad \alpha'_t = \alpha_t + n_{td},$$

где C — нормировочная константа, не зависящая от θ_d .

Оценим случайную величину θ_{td} её математическим ожиданием (3.1) по апостериорному распределению:

$$p(t|d, X_d, \alpha) = \int p(t|d) p(\theta_d|X_d, \alpha) d\theta_d = \int \theta_{td} \text{Dir}(\theta_d, \alpha') d\theta_d = \frac{n_{td} + \alpha_t}{n_d + \alpha_0}. \quad (3.2)$$

Заменяв величину n_{td} её оценкой \hat{n}_{td} , получим сглаженную байесовскую оценку параметра θ_{td} для EM-алгоритма, альтернативную оценке максимума правдоподобия (2.3):

$$\theta_{td} = \frac{\hat{n}_{td} + \alpha_t}{\hat{n}_d + \alpha_0}. \quad (3.3)$$

Аналогично выводится сглаженная байесовская оценка и для φ_{wt} , альтернативная (2.2):

$$\varphi_{wt} = \frac{\hat{n}_{wt} + \beta_w}{\hat{n}_t + \beta_0}. \quad (3.4)$$

Замена в обобщённом EM-алгоритме частотных оценок условных вероятностей (2.2) и (2.3) сглаженными оценками (3.4) и (3.3) трансформирует PLSA в LDA. Более строгое обоснование EM-подобных алгоритмов приводится в [30, 38] для метода сэмплирования Гиббса и в [32] для метода вариационной байесовской аппроксимации.

В [5] показано, что эти и другие известные алгоритмы обучения LDA являются вариантами EM-алгоритма и отличаются, главным образом, формулой сглаживания частотных оценок вероятностей. Оптимизация гиперпараметров α и β , предложенная в [35, 36], ещё сильнее нивелирует различия между моделями. Согласно экспериментам на 7 текстовых коллекциях [5], более эффективным по качеству и по времени является алгоритм *свёрнутой вариационной байесовской аппроксимации* CVB0 (collapsed variational Bayes). В нашей нотации ему наиболее близок LDA-GEM.

§3.2 Сэмплирование Гиббса

В задачах статистического оценивания часто возникает ситуация, когда вычисление или хранение некоторой функции распределении слишком ресурсоёмко, в то же время, генерация случайной выборки из этого распределения не вызывает затруднений. Сэмплированием Гиббса (Gibbs sampling, GS) называется общий метод решения

Алгоритм 3.1. LDA-GS: сэмплирование Гиббса для тематической модели LDA.

Вход: коллекция D , число тем $|T|$, начальные Θ , Φ , гиперпараметры α , β ;

Выход: распределения Θ и Φ ;

- 1 обнулить \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t для всех $d \in D$, $w \in W$, $t \in T$;;
 - 2 **повторять**
 - 3 **для всех** $d \in D$, $w \in d$, $i = 1, \dots, n_{dw}$
 - 4 **если** не первая итерация **то**
 - 5 $t := t_{dwi}$; уменьшить \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t на 1;
 - 6 сэмплировать тему t_{dwi} из $p(t | d, w) \propto (\hat{n}_{dt} + \alpha_t)(\hat{n}_{wt} + \beta_w)/(\hat{n}_t + \beta_0)$;
 - 7 $t := t_{dwi}$; увеличить \hat{n}_{wt} , \hat{n}_{dt} , \hat{n}_t на 1;
 - 8 **пока** Θ и Φ не сойдутся;
 - 9 $\varphi_{wt} = (\hat{n}_{wt} + \beta_w)/(\hat{n}_t + \beta_0)$ для всех $t \in T$, $w \in W$;
 - 10 $\theta_{td} := (\hat{n}_{dt} + \alpha_t)/(n_d + \alpha_0)$ для всех $d \in D$, $t \in T$;
-

таких задач, основанный на замене исходного распределения эмпирическим, вычисленным по выборке, сэмплированной из данного распределения.

Применение GS к тематической модели LDA предложено в [30]. Строгий вывод формул LDA-GS приводится в отчёте [38]. LDA-GS (Алгоритм 3.1) имеет несколько отличий от PLSA-SEM — стохастического варианта Алгоритма 2.2, но только одно из них оказывается существенным с точки зрения качества модели.

1. В LDA-GS жёстко фиксируется число сэмплирований тем $s = n_{dw}$ для каждой пары (d, w) . Однако *гипотеза разреженности* предполагает, что термин w в документе d связан с небольшим числом тем. В наших экспериментах $s = 5$ тем оказалось достаточно. В некоторых задачах достаточно и одной темы, в других одной темы мало, см. рис. 2. Эвристика *экономного сэмплирования* повышает эффективность алгоритма как по скорости, так и по памяти, но ухудшая качество модели.

Вставить график

2. В LDA-GS параметры φ_{wt} и θ_{td} обновляются предельно часто — после обработки каждого вхождения термина w в документ d . Эксперименты показывают, что достаточно делать обновления после каждой пары (d, w) , это не влияет на качество модели.

ToDo¹⁶

Вставить график

3. В LDA-GS перед сэмплированием счётчики уменьшаются на единицу (шаг 5). Тем самым в оценке распределений не учитывается i -е вхождение термина w в документ d , для которого сэмплируется тема t_{dwi} . Эта особенность алгоритма следует из теории [38]. Однако эксперименты показывают, что она не влияет на качество модели. Можно одновременно уменьшать счётчики для старой темы и увеличивать для новой, как в Алгоритме 2.2.

ToDo¹⁷

4. Единственным существенным различием, влияющим на качество модели, является применение в LDA байесовской регуляризации, приводящей к сглаживанию частотных оценок условных вероятностей.

Таким образом, LDA-GS существенно отличается от PLSA-GEM только тремя эвристиками: частотой обновления параметров, сэмплированием и сглаживанием. Эти эвристики не связаны друг с другом и могут применяться в любых сочетаниях.

§3.3 Оптимизация гиперпараметров

В первых работах по LDA [6] и сэмплированию Гиббса [30], а также в последовавших за ними исследованиях использовались симметричные распределения Дирихле с гиперпараметрами $\alpha = (a, \dots, a)$ и $\beta = (b, \dots, b)$. Скалярные гиперпараметры a и b либо фиксировались (одна из стандартных рекомендаций: $a = 50/|T|$, $b = 0.01$), либо настраивались путём перебора по сетке значений.

Позже были предложены численные методы оптимизации гиперпараметров, их обзор и сравнение приводится в диссертации [35]. Большинство методов оптимизации гиперпараметров основаны на максимизации *обоснованности* (evidence) модели, определяемой по всей коллекции $X = (X_d)_{d \in D}$:

$$\begin{aligned} P(X|\alpha) &= \int P(X|\theta)p(\theta|\alpha) d\theta = \int \prod_{d \in D} P(X_d|\theta_d) \text{Dir}(\theta_d; \alpha) d\theta = \\ &= \prod_{d \in D} \frac{\Gamma(\alpha_0)}{\Gamma(n_d + \alpha_0)} \prod_{t \in T} \frac{\Gamma(n_{td} + \alpha_t)}{\Gamma(\alpha_t)} \rightarrow \max_{\alpha} \end{aligned}$$

Метод неподвижной точки [24] — один из самых простых, но не самый лучший — представляет собой итерационный процесс:

$$\alpha_t := \alpha_t \frac{\sum_d \psi(n_{td} + \alpha_t) - \psi(\alpha_t)}{\sum_d \psi(n_d + \alpha_0) - \psi(\alpha_0)}, \quad t \in T,$$

где $\psi(z) = (\ln \Gamma(z))' = \Gamma'(z)/\Gamma(z)$ — дигамма-функция.

Этот или другой аналогичный итерационный процесс встраивается между проходами по коллекции, когда значения счётчиков n_{td} и n_d вычислены и фиксированы. Он выполняется намного быстрее одного прохода коллекции, поэтому оптимизацию гиперпараметров можно считать вычислительно эффективной.

Эксперименты показали, что оптимизация гиперпараметров существенно улучшает качество тематической модели [36]. Оказалось, что априорное распределение $\text{Dir}(\theta; \alpha)$ лучше брать несимметричным и оптимизировать вектор гиперпараметров $\alpha = (\alpha_1, \dots, \alpha_{|T|})$, а распределение $\text{Dir}(\varphi; \beta)$ лучше брать симметричным и оптимизировать скалярный гиперпараметр b , причём $0 < b \ll 1$.

§3.4 Позволяет ли сглаживание уменьшить переобучение

Эксперименты [6] показали, что LDA обеспечивает существенно меньшие значения контрольной перспексии, чем PLSA. По аналогии с задачами классификации и регрессии отсюда был сделан стандартный вывод, что модель PLSA имеет слишком много параметров θ_{td} , φ_{wt} , и при отсутствии ограничений на них возникает переобучение. Байесовская регуляризация должна сокращать эффективную размерность и уменьшать переобучение.

Однако возможна и другая интерпретация этих экспериментов. Сложность модели здесь не при чём — PLSA и LDA оценивают одни и те же матрицы параметров Φ и Θ . Оптимальные значения гиперпараметров α и β в LDA обычно близки к нулю. Поэтому оценки параметров φ_{wt} и θ_{td} в PLSA и в LDA могут заметно отличаться только для тем, редких в документе, и терминов, редких в теме. С одной стороны, они

не несут статистически значимой информации о тематике коллекции. С другой стороны, именно для редких терминов w тематическая модель предсказывает близкую к нулю вероятность $p(w | d)$. При появлении этих терминов в документах контрольная перплексия резко увеличивается. В ходе итераций PLSA оценки вероятности $p(w | d)$ редких терминов могут стремиться к нулю, что выглядит как переобучение, хотя по сути им не является. В LDA вероятности редких терминов никогда не стремятся к нулю благодаря сглаженным байесовским оценкам φ_{wt} и θ_{td} .

Возникает гипотеза, что контрольная перплексия у PLSA хуже, чем у LDA только из-за редких терминов, практически бесполезных для выявления тематики. Другими словами, кажущееся переобучение является побочным следствием гиперчувствительности перплексии к малым вероятностям.

Эксперименты [1] показали, что если из контрольных документов убрать новые термины, то сглаживание не даёт выигрыша, и перплексии PLSA и LDA практически совпадают, см. рис. 1. Этот результат согласуется с рядом недавних исследований [22, 39, 21], также подтверждающих, что для больших коллекций нет существенных различий в качестве моделей PLSA и LDA. Значимые отличия контрольной перплексии PLSA и LDA, наблюдавшиеся в ранних экспериментах [6], могут быть объяснены многочисленными различиями в реализации алгоритмов обучения. В экспериментах [1] реализации алгоритмов обучения PLSA и LDA отличались только сглаженными оценками в LDA.

Сглаживание создаёт ряд проблем: необходимо оптимизировать гиперпараметры, инициализировать β_w для новых терминов w и обеспечивать разреженность при том, что распределение Дирихле не позволяет обнулять вероятности θ_{td} и φ_{wt} .

Идея автоматического выделения терминов, бесполезных для тематической модели, приводит к робастным моделям, которые легко поддаются разреживанию и могут обходиться без байесовской регуляризации и сглаживания [1].

4 Робастные и разреженные тематические модели

Согласно вероятностной модели (1.2), каждый термин w в каждом документе d порождается некоторой темой t . Однако появление отдельных терминов может объясняться не только тематикой документа. Возможны, как минимум, ещё два альтернативных объяснения, условно называемых шумом и фоном.

Шум — это термины, специфичные для конкретного документа, либо редкие термины, относящиеся к темам, слабо представленным в данной коллекции. Тематическая модель даёт слишком низкие значения вероятности $p(w | d)$ для таких терминов, то есть не способна объяснить их появление в документах коллекции. Шумовые термины увеличивают перплексию и искажают тематическую модель.

Фон — это общеупотребительные слова, в частности, стоп-слова, не отброшенные на стадии предварительной обработки. Фоновые слова имеют значимые вероятности во многих темах, снижая релевантность тематического поиска.

§4.1 Робастная тематическая модель с шумом и фоном

Робастная вероятностная тематическая модель SWB (special words with background) представляет собой вероятностную смесь трёх компонент — тематиче-

ской, шумовой и фоновой [9]:

$$p(w | d) = \frac{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}{1 + \gamma + \varepsilon}; \quad Z_{dw} = \sum_{t \in T} \varphi_{wt}\theta_{td}. \quad (4.1)$$

где *шумовая компонента* $\pi_{dw} \equiv p_{\text{ш}}(w | d)$ — неизвестное распределение терминов в документе d , *фоновая компонента* $\pi_w \equiv p_{\text{ф}}(w)$ — неизвестное распределение терминов во всей коллекции. Априорные вероятности тематической, шумовой и фоновой компонент модели обозначим, соответственно, $q_{\text{т}} = \frac{1}{1+\gamma+\varepsilon}$, $q_{\text{ш}} = \frac{\gamma}{1+\gamma+\varepsilon}$, $q_{\text{ф}} = \frac{\varepsilon}{1+\gamma+\varepsilon}$, где γ и ε — неотрицательные параметры.

Суть робастной модели в том, что если тематическая компонента Z_{dw} плохо объясняет избыточную частоту n_{dw} некоторого термина w в некотором документе d , то она может быть объяснена альтернативным образом либо шумовой компонентной π_{dw} , либо фоновой π_w .

Требуется найти значения вероятностей $\varphi_{wt}, \theta_{td}, \pi_{dw}, \pi_w$, при которых логарифм правдоподобия достигает максимума:

$$L(D; \Theta, \Phi, \Pi) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \frac{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}{1 + \gamma + \varepsilon} \rightarrow \max_{\Theta, \Phi, \Pi}, \quad (4.2)$$

при ограничениях неотрицательности $\pi_{dw} \geq 0$, $\pi_w \geq 0$ и нормировки

$$\sum_{w \in W} \varphi_{wt} = 1, \quad \sum_{t \in T} \theta_{td} = 1, \quad \sum_{w \in d} \pi_{dw} = 1, \quad \sum_{w \in W} \pi_w = 1.$$

Чтобы получить приближённое решение М-шага, запишем лагранжиан данной задачи при ограничениях нормировки и неотрицательности π_{dw}, π_w , проигнорировав ограничения неотрицательности θ_{td} и φ_{wt} , которые будут выполнены автоматически.

$$\begin{aligned} \mathcal{L}(D; \Theta, \Phi, \Pi) &= \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \frac{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}{1 + \gamma + \varepsilon} + \sum_{d \in D} \sum_{w \in d} \kappa_{dw} \pi_{dw} + \sum_{w \in W} \kappa'_w \pi_w - \\ &- \sum_{t \in T} \lambda_t \left(\sum_{w \in W} \varphi_{wt} - 1 \right) - \sum_{d \in D} \mu_d \left(\sum_{t \in T} \theta_{td} - 1 \right) - \sum_{d \in D} \nu_d \left(\sum_{w \in d} \pi_{dw} - 1 \right) - \nu' \left(\sum_{w \in W} \pi_w - 1 \right). \end{aligned}$$

Двойственные переменные κ_{dw} , соответствующие ограничениям $\pi_{dw} \geq 0$, должны быть неотрицательны и удовлетворять условиям дополняющей нежёсткости

$$\kappa_{dw} \pi_{dw} = 0, \quad d \in D, \quad w \in d.$$

Аналогично, для двойственных переменных κ'_w , соответствующих $\pi_w \geq 0$:

$$\kappa'_w \pi_w = 0, \quad w \in W.$$

По аналогии со стандартным EM-алгоритмом, на E-шаге для каждой пары (d, w) вычисляются по формуле Байеса условные вероятности тем $H_{dwt} = p(t | d, w)$:

$$H_{dwt} = \frac{\varphi_{wt}\theta_{td}}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}, \quad t \in T, \quad (4.3)$$

а также условные вероятности того, что термин w является шумом H_{dw} и фоном H'_{dw} :

$$H_{dw} = \frac{\gamma\pi_{dw}}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}; \quad H'_{dw} = \frac{\varepsilon\pi_w}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w}. \quad (4.4)$$

Продифференцировав лагранжиан по переменным θ_{td} и φ_{wt} и приравняв нулю производные, получим прежние формулы для φ_{wt} (2.2) и θ_{td} (2.3), с единственным отличием, что теперь H_{dwt} вычисляются по новой формуле (4.3).

Продифференцируем лагранжиан по π_{dw} и приравняем нулю производную:

$$\nu_d = \frac{n_{dw}\gamma}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} + \kappa_{dw}. \quad (4.5)$$

Домножим обе части этого равенства на π_{dw} , просуммируем по всем терминам $w \in W$, применим условие нормировки вероятностей π_{dw} в левой части и условие дополняющей нежесткости в правой части. Получим выражение двойственной переменной ν_d через все основные переменные:

$$\nu_d = \sum_{w \in d} n_{dw} \frac{\gamma\pi_{dw}}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} = \sum_{w \in d} n_{dw} H_{dw}. \quad (4.6)$$

Поскольку H_{dw} есть апостериорная вероятность того, что термин w в документе d является шумом, величина ν_d интерпретируется как оценка числа шумовых терминов в документе d .

Проделав аналогичные действия для фоновой компоненты, получим

$$\begin{aligned} \nu' &= \sum_{d \in D} n_{dw} \frac{\varepsilon}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} + \kappa'_w, \\ \nu' &= \sum_{d \in D} \sum_{w \in d} n_{dw} \frac{\varepsilon\pi_w}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} = \sum_{d \in D} \sum_{w \in d} n_{dw} H'_{dw}, \end{aligned}$$

где ν' интерпретируется как оценка числа фоновых терминов во всей коллекции.

Мультипликативный М-шаг. Домножим обе части (4.5) на π_{dw} , но не будем суммировать по w . Получим формулу М-шага для шумовой компоненты:

$$\pi_{dw} = \frac{1}{\nu_d} n_{dw} \frac{\gamma\pi_{dw}}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} = \frac{n_{dw} H_{dw}}{\sum_{w' \in d} n_{dw'} H_{dw'}}.$$

Аналогично получается формула М-шага для фоновой компоненты:

$$\pi_w = \frac{1}{\nu'} n_{dw} \frac{\varepsilon\pi_w}{Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w} = \frac{\sum_{d \in D} n_{dw} H'_{dw}}{\sum_{d \in D} \sum_{w' \in d} n_{dw'} H'_{dw'}}.$$

Неотрицательность решения π_{dw} , π_w гарантируется, коль скоро начальные приближения π_{dw} , π_w неотрицательны. Мультипликативный М-шаг приводит к аналогичной проблеме разреженности для переменных π_{dw} и π_w , что и для переменных φ_{wt} и θ_{td} . Если в начальном приближении значение π_{dw} или π_w равно нулю, то оно сохранится и далее на протяжении итераций. Если в начальном приближении π_{dw} или π_w не равно нулю, то оно так и останется ненулевым.

Аддитивный М-шаг решает проблему разреживания шумовой компоненты [1]. Перепишем (4.5) в другом виде:

$$n_{dw}\gamma = (\nu_d - \kappa_{dw})(Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w).$$

Согласно условиям дополняющей нежёсткости, хотя бы одна из двух неотрицательных переменных κ_{dw} , π_{dw} должна быть равна нулю. Отсюда следует, что если $n_{dw}\gamma < \nu_d(Z_{dw} + \varepsilon\pi_w)$, то $\pi_{dw} = 0$ и $\kappa_{dw} > 0$. Если же имеет место противоположное неравенство, то $\kappa_{dw} = 0$ и π_{dw} находится из уравнения $n_{dw}\gamma = \nu_d(Z_{dw} + \gamma\pi_{dw} + \varepsilon\pi_w)$. Объединяя оба эти случая, получаем итоговое выражение для π_{dw} :

$$\pi_{dw} = \left(\frac{n_{dw}}{\nu_d} - \frac{Z_{dw} + \varepsilon\pi_w}{\gamma} \right)_+. \quad (4.7)$$

Таким образом, если термин w в документе d встречается существенно чаще, чем предсказывают тематическая и фоновая компоненты модели, то его появление объясняется особенностями данного документа, и тогда $\pi_{dw} > 0$.

Аддитивный М-шаг, в отличие от мультипликативного, приводит к автоматическому выбору структуры разреженности матрицы $(\pi_{dw})_{D \times W}$.

Эксперименты показывают, что аддитивный М-шаг даёт небольшое преимущество по перплексии. ToDo¹⁸

Робастная модификация PLSA-ROEM онлайнowego PLSA-OEM показана в Алгоритме 4.1. Главное отличие от обычного PLSA в том, что теперь n_{dw} вхождений термина w в документ d распределяются не только между темами $t \in T$, но также между шумовой и фоновой компонентами, пропорционально вероятностям

$$\tilde{H}_{dw} = \left(\frac{1}{Z}\varphi_{wt}\theta_{td}, t \in T; \frac{1}{Z}\gamma\pi_{dw}; \frac{1}{Z}\varepsilon\pi_w \right),$$

где Z — нормирующий множитель.

Возможны различные варианты алгоритма PLSA-ROEM: только с шумовой компонентой ($\varepsilon = 0$), только с фоновой компонентой ($\gamma = 0$), с аддитивным и мультипликативным М-шагом. В Алгоритме 4.1 показан вариант с шумом и фоном, аддитивным М-шагом, без сэмплирования и без сглаживания.

Сглаживание вводится в Алгоритм 4.1 заменой частотных оценок (2.2)–(2.3) параметров φ_{wt} , θ_{td} на шагах 17, 12 байесовскими оценками (3.3)–(3.4).

Сэмплирование вводится заменой распределения \tilde{H}_{dw} его эмпирической оценкой, аналогичной (2.5).

Результаты экспериментов: робастная модель менее чувствительна к выбору параметра экономного сэмплирования s . ToDo¹⁹

Два варианта сэмплирования для каждого (d, w) :

- (1) пропорциональное распределение вероятности темы–шум–фон, темы сэмплируются;
 - (2) сэмплирование из всего распределения \tilde{H}_{dw} .
- ToDo²⁰

О невозможности оптимизации априорных вероятностей шума и фона. Приравняв нулю производные лагранжиана по γ и ε , нетрудно получить формулы для обновления γ и ε . Однако эксперименты показывают, что с итерациями $\gamma \rightarrow \infty$, $\varepsilon \rightarrow 0$, что приводит к полному вырождению тематической модели в простейшую униграммную модель. Поэтому параметры γ и ε необходимо фиксировать.

Возможна оптимизация с помощью непараметрического байесовского вывода. Правда, тогда появится ещё один параметр, меняя который можно всё убрать в шум... ToDo²¹

Алгоритм 4.1. PLSA-ROEM: робастный онлайнный EM-алгоритм.

Вход: коллекция документов D , число тем $|T|$;

Выход: распределения Θ , Φ , Π ;

```

1 инициализировать  $\varphi_{wt}$ ,  $\pi_w$  для всех  $w \in W$ ,  $t \in T$ ;
2  $\hat{n}_{wt} := 0$ ,  $\hat{n}_t := 0$ ,  $\hat{n}'_w := 0$ ,  $\hat{n}' := 0$  для всех  $w \in W$ ,  $t \in T$ ;
3 для всех пакетов  $D_j$ ,  $j = 1, \dots, J$ 
4   повторять
5      $\tilde{n}_{wt} := 0$ ,  $\tilde{n}_t := 0$ ,  $\tilde{n}'_w := 0$ ,  $\tilde{n}' := 0$  для всех  $w \in W$ ,  $t \in T$ ;
6     для всех  $d \in D_j$ 
7       инициализировать  $\theta_{td}$  для всех  $t \in T$ ;
8       повторять
9          $Z_w := \sum_{t \in T} \varphi_{wt} \theta_{td} + \gamma \pi_{dw} + \varepsilon \pi_w$  для всех  $w \in d$ ;
10         $\nu_d := \sum_{w \in d} n_{dw} \gamma \pi_{dw} / Z_w$ ;
11         $\tilde{n}_t := \sum_{w \in d} n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $t \in T$ ;
12         $\theta_{td} := \tilde{n}_t / \sum_{s \in T} \tilde{n}_s$  для всех  $t \in T$ ;
13         $\pi_{dw} := (\pi_{dw} + n_{dw} / \nu_d - Z_w / \gamma)_+$  для всех  $w \in d$ ;
14        пока  $\theta_{td}$  и  $\pi_{dw}$  для данного  $d$  не сойдутся;
15        увеличить  $\tilde{n}_{wt}$ ,  $\tilde{n}_t$  на  $n_{dw} \varphi_{wt} \theta_{td} / Z_w$  для всех  $w \in d$ ,  $t \in T$ ;
16        увеличить  $\tilde{n}'_w$ ,  $\tilde{n}'$  на  $n_{dw} \varepsilon \pi_w / Z_w$  для всех  $w \in d$ ;
17       $\varphi_{wt} := \frac{\rho_j \hat{n}_{wt} + \tilde{n}_{wt}}{\rho_j \hat{n}_t + \tilde{n}_t}$  для всех  $w \in W$ ,  $t \in T$ ;
18       $\pi_w := \frac{\rho_j \hat{n}'_w + \tilde{n}'_w}{\rho_j \hat{n}' + \tilde{n}'}$  для всех  $w \in W$ ;
19    пока  $\Phi$  не сойдутся;
20     $\hat{n}_{wt} := \rho_j \hat{n}_{wt} + \tilde{n}_{wt}$  для всех  $w \in W$ ,  $t \in T$ ;
21     $\hat{n}_t := \rho_j \hat{n}_t + \tilde{n}_t$  для всех  $t \in T$ ;
22     $\hat{n}'_w := \rho_j \hat{n}'_w + \tilde{n}'_w$  для всех  $w \in W$ ;
23     $\hat{n}' := \rho_j \hat{n}' + \tilde{n}'$ ;

```

§4.2 Принудительное разреживание

Гипотеза разреженности предполагает, что в дискретных распределениях $p(w|t) = \varphi_{wt}$, $p(t|d) = \theta_{td}$, $p(t|d, w) = H_{dwt}$ подавляющее большинство вероятностей равны нулю или очень близки к нулю. Алгоритмы, в которых нулевые значения не хранятся, намного эффективнее по памяти и по скорости. Поэтому для больших коллекций разреженность должна учитываться обязательно.

Модель PLSA не оптимизирует структуру разреженности распределений и требует задавать её через начальное приближение. Отдельные значения θ_{td} и φ_{wt} могут в ходе итераций стремиться к нулю, но, как правило, их доля недостаточна для получения выигрыша в производительности.

Модель LDA также не является разреженной — сглаживание частотных оценок вероятностей приводит к тому, что матрицы Φ и Θ не содержат нулевых значений.

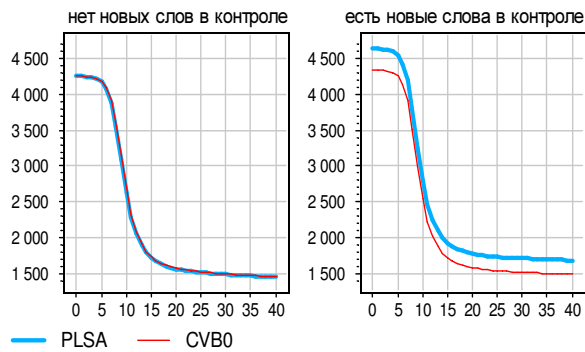


Рис. 1. Сглаживание даёт преимущество только когда в контроле есть новые термины (метод CVB0 — это PLSA-GEM со сглаживанием но без сэмплирования).

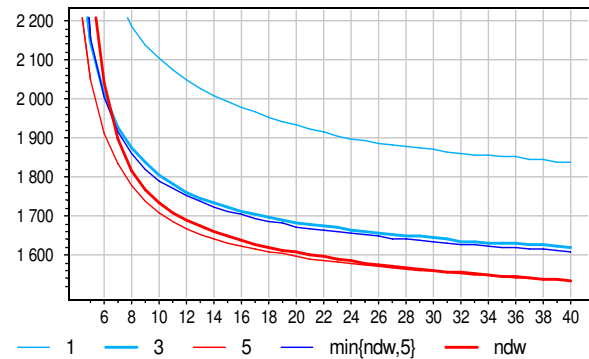


Рис. 2. При экономном сэмплировании пяти тем для каждой пары (d, w) перплексия не хуже, чем при сэмплировании n_{dw} тем. Но одной или трёх тем не достаточно.

Эта проблема имеет много решений, например в [13] предлагается хранить не сами значения θ_{td} и φ_{wt} , а только их разности с фоновыми распределениями.

Принудительное разреживание вводится в любой из описанных выше EM-подобных алгоритмов. По окончании цикла по всем документам $d \in D$ в каждом из $|T|$ распределений $\varphi_{wt} = p(w|t)$ обнуляется небольшая доля r_φ наименьших значений вероятностей. Аналогично, в каждом из $|D|$ распределений $\theta_{td} = p(t|d)$ обнуляется небольшая доля r_θ наименьших значений. После обнуления производится перенормировка распределений. Разреживания начинаются с некоторой итерации i_0 , чтобы к этому моменту в распределениях правильно выделились малые вероятности. Кроме того, разреживания имеет смысл делать не на каждой итерации, а хотя бы через одну, чтобы восстановить адекватность модели. В наших экспериментах разреживание выполнялось на итерациях с номерами вида $i = i_0 + kd$, $k = 1, 2, \dots$ при $r_\varphi = 0.05$, $r_\theta = 0.05$, $i_0 = 10$, $d = 2$.

Принудительное разреживание может увеличивать перплексию моделей PLSA и LDA. Кроме того, при сильном разреживании в PLSA может происходить обнуление модели $p(w|d) = 0$ при $n_{dw} > 0$, и тогда перплексия уходит в бесконечность.

Робастные модели допускают принудительное разреживание без ухудшения перплексии и одновременно исключают ситуацию бесконечной перплексии, так как нулевое значение Z_{dw} компенсируется ненулевым значением шума π_{dw} или фона π_w . Чем больше γ и ε , тем сильнее можно разредить тематическую компоненту модели. В экспериментах разреженность матриц Θ и Φ в робастном PLSA достигала более 90% без потери качества модели или даже с незначительным улучшением, рис. 4.

Эксперименты: сравнение разреживания для PLSA, LDA и робастных.

ToDo²²

Эксперименты с выбором стратегии разреживания.

ToDo²³

Эксперименты на реальных данных Эксперименты производились на двух коллекциях.

Коллекция RuDis содержала $|D| = 2000$ авторефератов диссертаций на русском языке; суммарная длина $n \approx 8.7 \cdot 10^6$, объём словаря $|W| \approx 3 \cdot 10^4$. Контрольная коллекция D' состояла из 200 авторефератов. Предварительно производилась лемматизация и отбрасывались стоп-слова.

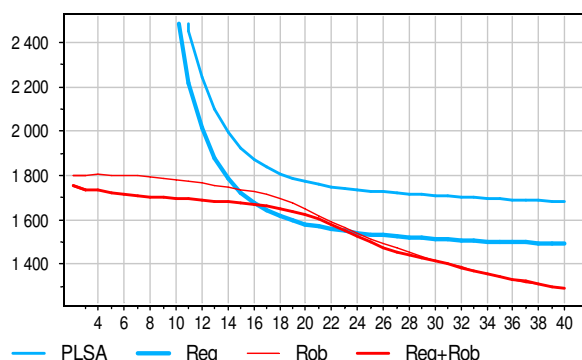


Рис. 3. Робастность сильнее уменьшает перплексию PLSA, чем сглаживание. Сглаживание не улучшает робастную модель.

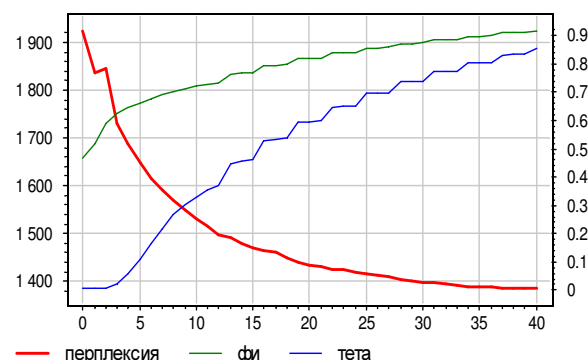


Рис. 4. В процессе разреживания доля нулевых φ_{wt} и θ_{td} (отложена по правой оси) увеличивается при монотонном уменьшении перплексии.

Коллекция NIPS содержала $|D| = 1566$ текстов статей научной конференции Neural Information Processing Systems на английском языке; суммарная длина $n \approx 2.3 \cdot 10^6$, объём словаря $|W| \approx 1.3 \cdot 10^4$. Контрольная коллекция D' состояла из 174 документов. Предварительно производился стемминг и отбрасывались стоп-слова.

Качество модели оценивалось *перплексией* контрольной коллекции D' документов, не включённых в обучающую коллекцию. Каждый контрольный документ d случайным образом разделялся на две половины, d' и d'' . Параметры θ_{td} и ν_d оценивались по d' . Параметры φ_{wt} и π_w оценивались по обучающей выборке D . Параметры π_{dw} оценивались для каждой пары (d, w) согласно (4.7). Перплексия вычислялась по вторым половинам d'' контрольных документов.

На рис. 1–4 показаны зависимости перплексии от числа итераций (одна итерация — один проход по коллекции). Число итераций 40; число тем $|T| = 100$; параметры сглаживания $\alpha_t = 0.5$, $\beta_w = 0.01$; параметры робастности $\gamma = 0.3$, $\varepsilon = 0.1$.

Проверить, действительно ли $\pi_w > 0$ для стоп-слов; $\pi_{dw} > 0$ для слов, не типичных ни для одной из тем. Просмотреть эти слова и сравнить их со словами из тем. Показать одни и те же темы как ранжированные списки терминов в неробастной модели, робастной модели, робастной модели после разреживания. ToDo²⁴

Если добавить фоновую компоненту, то приведёт ли это к улучшению НОР и степени разреженности? Позволяет ли фоновая компонента избавиться от предварительного отбрасывания стоп-слов? ToDo²⁵

5 Иерархические тематические модели

Для больших коллекций текстовых документов естественно строить иерархии вложенных друг в друга тем (называемых также категориями или рубриками), чтобы упростить поиск документов. Иерархия — это общепринятый способ структуризации знаний. Однако разделение тем на более узкие подтемы субъективно, неоднозначно и часто вызывает споры среди специалистов.

В статье [42] приводится обзор иерархических тематических моделей и отмечается, что оптимизация структуры иерархии по коллекции документов является

открытой проблемой; более того, разработка объективной количественной оценки качества иерархии — также открытая проблема.

Многие иерархические модели имеют те или иные неестественные ограничения: либо фиксируется число уровней, либо фиксируется число подтем в каждой теме или на каждом уровне, либо документ не может относиться к темам из различных ветвей дерева, либо темы не могут иметь общую подтему, либо темам во внутренних узлах не сопоставляется распределение на множестве терминов.

§5.1 Определение тематического дерева

Гипотеза о существовании тематического дерева. Рассмотрим дерево с множеством вершин V и корнем $t_0 \in V$. Вершины дерева соответствуют темам. Каждой теме $t \in V$ соответствует множество её подтем — дочерних вершин в дереве $S_t \subset V$. Каждое ребро дерева соответствует паре «тема–подтема» (t, s) , $s \in S_t$. Если $S_t = \emptyset$, то тема t называется *терминальной* или *листом* тематического дерева. Для каждой вершины t в дереве V существует только одна родительская вершина, следовательно, только один путь (t_0, \dots, t) от корня дерева t_0 до темы t .

Ранее мы предполагали, что каждое вхождение термина w в документ d связано только с одной темой t . Теперь примем за аксиому другие предположения:

- 1) если пара (d, w) связана с темой t , то она связана и со всеми темами выше вершины t на пути до корня t_0 ;
- 2) если пара (d, w) не связана с темой t , то она не связана и со всеми подтемами в поддереве ниже вершины t .

Этих двух предположений, совершенно не вероятностного характера, достаточно, чтобы построить иерархическую вероятностную тематическую модель.

Вероятностная интерпретация отношения «тема–подтема». Каждому ребру тематического дерева (t, s) соответствует условная вероятность $p(s | t)$ того, что термин документа, связанный с темой t , связан также с подтемой $s \in S_t$:

$$p(s | t) = \frac{p(t, s)}{p(t)} = \frac{p(s)}{p(t)}. \quad (5.1)$$

Если рассматривать коллекцию документов как выборку троек (d, w, t) , то частотной оценкой этой условной вероятности будет $\hat{p}(s | t) = n_s/n_t$ — доля троек, связанных с подтемой s , среди всех троек, связанных с темой t .

Условные вероятности подтем удовлетворяют ограничениям нормировки, которые, в силу (5.1), допускают две эквивалентные записи:

$$\sum_{s \in S_t} p(s | t) = 1, \quad \sum_{s \in S_t} p(s) = p(t), \quad t \in V. \quad (5.2)$$

Обозначим через T множество тем, соответствующих терминальным вершинам дерева V . Условие нормировки

$$\sum_{t \in T} p(t) = 1. \quad (5.3)$$

выполняется именно для этого множества, а не для всего множества тем в дереве V . Из (5.2) следует, что условие нормировки останется справедливым, если заменить

любое из множеств S_t его родительской темой t , а также если делать такие замены многократно в произвольном порядке. В частности, для корневой темы $p(t_0) = 1$.

При разделении темы t на подтемы $s \in S_t$ условные распределения для подтем $\varphi_{ws} = p(w | s)$ и $\theta_{sd} = p(s | d)$ должны удовлетворять требованиям нормировки

$$\sum_{w \in W} p(w | s) = 1, \quad s \in S_t; \quad \sum_{s \in S_t} p(s | d) = p(t | d), \quad d \in D. \quad (5.4)$$

Распределения $p(s | w) = p(w | s) \frac{p(s)}{p(w)}$ и $p(d | s) = p(s | d) \frac{p(d)}{p(s)}$ также должны быть нормированы, откуда следуют ещё две серии тождеств:

$$\sum_{s \in S_t} p(w | s) p(s) = p(w | t) p(t), \quad w \in W; \quad \sum_{d \in D} p(s | d) p(d) = p(s), \quad s \in S_t. \quad (5.5)$$

Документы во внутренних вершинах. В некоторых приложениях важно, чтобы документы и термины могли относиться не только к терминальным вершинам, но и к любым внутренним вершинам тематического дерева. В частности, это могут быть документы, относящиеся сразу к нескольким подтемам, либо новые документы, которые пока не выделились в отдельную подтему.

Для каждой внутренней вершины $t \in V \setminus T$ создаётся выделенная терминальная вершина — подтема $s_0 \in S_t$. Если документ или термин попадает в s_0 , то считается, что он остался в теме t . В терминах кластеризации выделенная подтема s_0 — это специальный «фоновый» кластер, к которому относится всё, что не удалось с уверенностью отнести к другим кластерам — подтемам темы t .

К выделенной подтеме s_0 естественно предъявлять требование минимизации числа документов и описывать её тем же распределением, что и родительскую тему t .

§5.2 Иерархическая модель для категоризации текстов

Рассмотрим случай, когда структура тематического дерева $\{S_t : t \in V\}$ фиксирована. Чтобы каждая тема $t \in T$ имела интерпретацию, к ней привязывается множество документов $D_t \subset D$ и множество терминов $W_t \subset W$. Одно из этих множеств может быть пустым. Таким образом, ставится задача *частичного обучения* (semi-supervised learning) иерархической тематической модели.

Распределение $\theta_{td} = p(t | d)$, полученное в результате тематического моделирования, непосредственно решает задачу категоризации — документ d относится к тем темам (категориям), для которых вероятность θ_{td} превышает заданный порог. Для решения задач категоризации лучше подходят разреженные модели, в которых малые вероятности θ_{td} обнуляются в процессе построения модели. Обнуление или игнорирование малых вероятностей в уже построенной модели может приводить к менее адекватным результатам.

Для категоризации текстов часто применяется другой подход: для каждой пары тема–подтема строится классификатор на два класса [29]. Каждый классификатор обучается по выборке документов, относящихся к родительской теме, что требует больших затрат времени и памяти. Иерархическая тематическая модель, очевидно, является более естественным инструментом категоризации.

Иерархический Алгоритм 5.1 основан на онлайн-овом Алгоритме 2.4. Основное отличие PLSA-НОЕМ в том, что для вычисления распределения θ_{td} тем t в документе d производится спуск по дереву от корня к терминальным вершинам.

Модификация формул М-шага для иерархической модели. Пусть $T \subset V$ — подмножество вершин дерева, удовлетворяющее условию нормировки (5.3), и для всех тем $t \in T$ известны значения параметров φ_{wt} , θ_{td} . Сначала T состоит из единственной корневой вершины t_0 , для которой $\varphi_{wt_0} = p(w)$ и $\theta_{t_0d} = 1$. Спуск по дереву — это итерационный процесс, на каждом шаге которого выбирается некоторое подмножество тем $R \subseteq T$, и для каждой вершины $s \in S$ из множества всех их подтем $S = \bigcup_{t \in R} S_t$, вычисляются значения параметров φ_{ws} , θ_{sd} . После этого множество T заменяется на $(T \setminus R) \cup S$ и начинается следующий шаг. Спуск продолжается, пока T не совпадёт с множеством терминальных вершин дерева.

Рассмотрим вероятностную модель (1.2) при ограничениях нормировки (5.4). Параметры φ_{wt} и θ_{td} для всех тем $t \in T \setminus R$ будем считать фиксированными. Обозначим через σ_{dw} фиксированную часть вероятностной тематической модели:

$$\sigma_{dw} = \sum_{t \in T \setminus R} \varphi_{wt} \theta_{td}, \quad d \in D, w \in d.$$

Задача оценивания параметров φ_{ws} , θ_{sd} , $s \in S$ сводится к максимизации логарифма правдоподобия, аналогично задаче (1.6), но оптимизируется только часть параметров $\Phi_S = (\varphi_{ws})_{W \times S}$ и $\Theta_S = (\theta_{sd})_{S \times D}$, связанных с темами из S :

$$\begin{aligned} L(D; \Theta_S, \Phi_S) &= \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \left(\sigma_{dw} + \sum_{s \in S} \varphi_{ws} \theta_{sd} \right) \rightarrow \max_{\Phi_S, \Theta_S}; \\ &\sum_{w \in W} \varphi_{ws} = 1, \quad s \in S; \\ &\sum_{s \in S_t} \theta_{sd} = \theta_{td}, \quad t \in R, d \in D. \end{aligned}$$

Как обычно, нужно записать лагранжиан, приравнять нулю его производные по переменным φ_{ws} и θ_{sd} , из полученных уравнений исключить двойственные переменные и выразить φ_{ws} и θ_{sd} через H_{dws} :

$$\begin{aligned} H_{dws} &= \frac{\varphi_{ws} \theta_{sd}}{\sigma_{dw} + \sum_{s' \in S} \varphi_{ws'} \theta_{s'd}}, \quad d \in D, w \in d, s \in S; \\ \varphi_{ws} &= \frac{\sum_{d \in D} n_{dw} H_{dws}}{\sum_{w' \in W} \sum_{d \in D} n_{dw'} H_{dw's}}, \quad w \in W, s \in S; \\ \theta_{sd} &= \theta_{td} \frac{\sum_{w \in d} n_{dw} H_{dws}}{\sum_{s' \in S_t} \sum_{w' \in d} n_{dw'} H_{dw's'}}, \quad d \in D, s \in S_t, t \in R; \end{aligned}$$

или, в более компактной записи с использованием счётчиков:

$$\varphi_{ws} = \frac{\hat{n}_{ws}}{\hat{n}_s}, \quad \hat{n}_s = \sum_{w \in W} \hat{n}_{ws}, \quad \hat{n}_{ws} = \sum_{d \in D} n_{dw} H_{dws}. \quad (5.6)$$

$$\theta_{sd} = \theta_{td} \frac{\hat{n}_{ds}}{\hat{n}_{dt}}, \quad \hat{n}_{dt} = \sum_{s \in S_t} \hat{n}_{ds}, \quad \hat{n}_{ds} = \sum_{w \in d} n_{dw} H_{dws}. \quad (5.7)$$

Таким образом, формулы М-шага и Е-шага для иерархического алгоритма лишь немногим отличаются от обычного PLSA-EM.

Алгоритм 5.1. PLSA-НОЕМ: иерархический онлайнный EM-алгоритм.

Вход: коллекция документов D ; параметры λ и μ ,
множество тем V и структура тематического дерева $\{S_t: t \in V\}$,
привязки терминов и документов к темам φ_{wt}^0 и θ_{td}^0 ;
Выход: распределения Θ и Φ ;

```

1 инициализировать  $\varphi_{wt}$  с учётом  $\varphi_{wt}^0$  для всех  $w \in W, t \in V$ ;
2 повторять
3    $\hat{n}_{wt} := 0; \hat{n}_t := 0$  для всех  $w \in W, t \in V$ ;
4   для всех  $d \in D$ 
5     инициализировать  $\theta_{td}$  с учётом  $\theta_{td}^0$  для всех  $t \in V$ ;
6      $T := \{t_0\}; R := \{t_0\}; \theta_{t_0d} = 1$ ;
7     пока множество подтем  $S := \bigcup_{t \in R} S_t$  не пусто
8        $\sigma_{dw} = \sum_{t \in T \setminus R} \varphi_{wt} \theta_{td}$  для всех  $w \in d$ ;
9       повторять
10         $Z_w := \sigma_{dw} + \sum_{s \in S} \varphi_{ws} \theta_{sd}$  для всех  $w \in d$ ;
11         $\theta_{sd} := \mu \theta_{sd}^0 \theta_{td} + (1 - \mu) \theta_{td} n_d^{-1} \sum_{w \in d} n_{dw} \varphi_{ws} \theta_{sd} / Z_w$  для всех  $s \in S_t$ ,
12         $t \in R$ ;
13        пока  $\theta_{sd}$  не сойдутся для всех  $s \in S$ ;
14        увеличить  $\hat{n}_{ws}, \hat{n}_s$  на  $n_{dw} \varphi_{ws} \theta_{sd} / Z_w$  для всех  $w \in d, s \in S$ ;
15         $T := (T \setminus R) \cup S; R := S$ ;
16    $\varphi_{wt} := \lambda \varphi_{wt}^0 + (1 - \lambda) \hat{n}_{wt} / \hat{n}_t$  для всех  $w \in W, t \in V$ ;
17 пока  $\Phi$  не сойдутся;

```

Частичное обучение. Привязки терминов и документов к темам задаются в виде начальных распределений φ_{wt}^0 и θ_{td}^0 , вычисляемых по формулам из §2.5. Они используются как при инициализации распределений, так и на каждом шаге EM-алгоритма. При инициализации они смешиваются с неразрезанными случайными распределениями. На каждом шаге они смешиваются с оценками (5.6) и (5.7). Благодаря тому, что текущие приближения φ_{wt} и θ_{td} немного притягиваются к начальным распределениям φ_{wt}^0 и θ_{td}^0 , они не уходят далеко от исходно заданных интерпретаций. Сила этого «притяжения» регулируется параметрами λ и μ .

Регуляризация Дирихле может быть добавлена в Алгоритм 5.1 обычным образом: частотные оценки условных вероятностей (5.6), (5.7) заменяются сглаженными:

$$\varphi_{ws} = \frac{\beta_w + \hat{n}_{ws}}{\sum_{w' \in W} (\beta_{w'} + \hat{n}_{w's})} = \frac{\beta_w + \hat{n}_{ws}}{\beta_0 + \hat{n}_s}. \quad (5.8)$$

$$\theta_{sd} = \theta_{td} \frac{\alpha_s + \hat{n}_{ds}}{\sum_{s' \in S_t} (\alpha_{s'} + \hat{n}_{ds'})} = \theta_{td} \frac{\alpha_s + \hat{n}_{ds}}{\alpha_t + \hat{n}_{dt}}. \quad (5.9)$$

Заметим, что при разделении темы t на множество подтем S_t расщепляются также и гиперпараметры распределения Дирихле $\text{Dir}(\theta_d; \alpha)$, а их сумма α_0 не меняется. Это следует из условий нормировки (5.4) и свойства (3.1):

$$\sum_{s \in S_t} \theta_{sd} = \theta_{td} \quad \Rightarrow \quad \sum_{s \in S_t} \mathbb{E} \theta_{sd} = \mathbb{E} \theta_{td} \quad \Rightarrow \quad \sum_{s \in S_t} \alpha_s = \alpha_t.$$

6 Критерии качества тематических моделей

§6.1 Критерии условной независимости

Гипотеза условной независимости $p(w | d, t) = p(w | t)$ чрезвычайно важна для вероятностных тематических моделей. Именно она обеспечивает переход к компактному представлению данных $F \approx \Phi\Theta$. Для её проверки не требуется выделять контрольную выборку, что является преимуществом данного типа критериев.

Оба распределения легко оцениваются в EM-алгоритме:

$$\hat{p}(w | d, t) = \frac{n_{dwt}}{\hat{n}_{dt}}, \quad t \in T, d \in D;$$

$$\hat{p}(w | t) = \frac{\hat{n}_{wt}}{\hat{n}_t}, \quad t \in T.$$

Рассмотрим статистические тесты, проверяющие нулевую гипотезу о том, что различия между этими распределениями незначимы, точнее, что выборка с эмпирическим распределением $\hat{p}(w | d, t)$ могла быть получена из генеральной совокупности с распределением $\hat{p}(w | t)$.

Число \hat{n}_{dt} интерпретируется как длина части документа d , связанной с темой t , а число n_{dwt} — как число вхождений термина w в документ d , связанных с темой t . Введём ожидаемое число вхождений термина w в документ d , связанных с темой t :

$$E_{dwt} = \hat{n}_{dt} \hat{p}(w | t).$$

Критерий χ^2 Пирсона основан на вычислении статистики хи-квадрат, которая является естественной мерой различия двух распределений:

$$X_{dt}^2 = \sum_{w \in W_{dt}} \frac{(E_{dwt} - n_{dwt})^2}{E_{dwt}} = \hat{n}_{dt} \sum_{w \in W_{dt}} \frac{(\hat{p}(w | t) - \hat{p}(w | d, t))^2}{\hat{p}(w | t)},$$

где $W_{dt} = \{w \in W : E_{dwt} > 0\}$. Если значение X_{dt}^2 превышает $(1 - \alpha)$ -квантиль распределения хи-квадрат $\chi_{k, 1-\alpha}^2$ с числом степеней свободы $k = |W_{dt}| - 1$, то нулевая гипотеза отвергается.

Условием применимости асимптотики χ_k^2 считается наличие достаточного числа наблюдений во всей выборке, $\hat{n}_{dt} \geq 50$, а также достаточного ожидаемого числа наблюдений каждого термина, $E_{dwt} \geq 5$. Второе требование в типичном случае не выполняется для большинства терминов w , так как распределение $\hat{p}(w | t)$, как правило, разрежено, более того, мощность словаря W_{dt} может превышать длину документа \hat{n}_{dt} . Таким образом, в нашем случае критерий Пирсона применять нельзя. Для случая разреженных распределений больше подходят статистики G^2 и D^2 .

Статистика G^2 определяется через дивергенцию Кульбака–Лейблера, и для неё также справедливо асимптотическое распределение χ_k^2 с тем же числом степеней свободы, но при менее жёстких требованиях к числу наблюдений:

$$G_{dt}^2 = 2 \sum_{w \in W_{dt}} n_{dwt} \ln \frac{n_{dwt}}{E_{dwt}}.$$

Статистика D^2 — это поправка к статистике X^2 , предложенная Зельтерманом в [43] специально для случая разреженных распределений:

$$D_{dt}^2 = \sum_{w \in W_{dt}} \frac{(E_{dwt} - n_{dwt})^2 - n_{dwt}}{E_{dwt}}.$$

Эта статистика имеет асимптотически нормальное распределение. Особенности её применения обсуждаются в [31, 18].

Семейство функций расстояния Кресси–Рида. Для сравнения эмпирической функции вероятности $\hat{p}(w)$, оцененной по выборке длины n , с истинной функцией вероятности $p(w)$ принято использовать функции расстояния, придающие больший вес малым вероятностям:

$$\text{KL}(\hat{p} \| p) = \sum_w \hat{p}(w) \ln \frac{\hat{p}(w)}{p(w)} \text{ — дивергенция Кульбака–Лейблера;} \quad (6.1)$$

$$X^2(\hat{p}, p) = \sum_w \frac{(p(w) - \hat{p}(w))^2}{p(w)} \text{ — ненормированная } \chi^2\text{-статистика;} \quad (6.2)$$

$$H^2(\hat{p}, p) = \sum_w \left(\sqrt{p(w)} - \sqrt{\hat{p}(w)} \right)^2 \text{ — расстояние Хеллингера.} \quad (6.3)$$

Эти и другие «разумные» функции расстояния обобщаются (с точностью до константного множителя) параметрическим семейством дивергенций Кресси–Рида [10, 27]:

$$\text{CR}_\lambda(\hat{p} : p) = \frac{2}{\lambda(\lambda + 1)} \sum_w \hat{p}(w) \left(\left(\frac{\hat{p}(w)}{p(w)} \right)^\lambda - 1 \right).$$

Перестановочный тест основан на использовании эмпирического распределения статистики, полученного путём сэмплирования большого числа выборок в условиях истинности нулевой гипотезы. Перестановочные тесты применяются в тех случаях, когда функция распределения статистики неизвестна или имеет слишком сложный вид или её известные асимптотики не достаточно точны.

Пусть S — одна из статистик X^2 , G^2 , D^2 . Зафиксируем тему t . Сгенерируем N независимых выборок терминов из распределения $\hat{p}(w | t)$. Для каждой из них вычислим эмпирическое распределение $\hat{p}(w)$ и значение статистики S . По выборке значений статистики $\{S_1, \dots, S_N\}$ построим эмпирическое распределение $\hat{F}_t(S)$ и найдём его $(1 - \alpha)$ -квантиль $\hat{F}_{t, 1-\alpha}$. Число N должно быть порядка 10^3 при $\alpha = 0.05$.

Обозначим через S_{dt} значение статистики S , вычисленное по распределению $\hat{p}(w | d, t)$ для заданных $t \in T$ и $d \in D$. Поскольку распределение $\hat{F}_t(S)$ построено

в условиях истинности нулевой гипотезы, неравенство $S_{dt} > \hat{F}_{t,1-\alpha}$ является критерием отклонения нулевой гипотезы для документа d на уровне значимости α .

Заметим, что квантиль $\hat{F}_{t,1-\alpha}$ достаточно вычислить один раз для каждой темы t и использовать для всех документов $d \in D$, что даёт значительную экономию времени. Однако при изменении распределения $\hat{p}(w | t)$ распределение $\hat{F}_t(S)$ и его квантиль придётся пересчитать заново.

Оценки средней несогласованности для документов и тем. Введём индикатор события «тема t не согласована в документе d при уровне значимости α »:

$$B_{dt}(\alpha) = [S_{dt} > \hat{F}_{t,1-\alpha}];$$

Определим *среднюю несогласованность* темы, документа и тематической модели в целом при уровне значимости α :

$$B_t(\alpha) = \sum_{d \in D} \frac{\hat{n}_{dt}}{\hat{n}_t} B_{dt}(\alpha) \quad \text{— средняя несогласованность темы } t;$$

$$B_d(\alpha) = \sum_{t \in T} \frac{\hat{n}_{dt}}{n_d} B_{dt}(\alpha) \quad \text{— средняя несогласованность документа } d;$$

$$B(\alpha) = \sum_{d \in D} \sum_{t \in T} \frac{\hat{n}_{dt}}{n} B_{dt}(\alpha) \quad \text{— средняя несогласованность модели.}$$

Это нормированные величины, принимающие значения из отрезка $[0, 1]$. Чем меньше средняя несогласованность, тем лучше модель описывает соответствующую тему t , документ d или всю коллекцию в целом.

Критерий условной независимости. В [23] предлагается ещё один критерий, оценивающий степень несоответствия темы $t \in T$ гипотезе условной независимости. Он основан на дивергенции Кульбака–Лейблера и может быть легко вычислен в EM-алгоритме на каждом проходе коллекции:

$$\text{KL}_t = \text{KL}(\hat{p}(d, w | t) \parallel \hat{p}(d | t) \hat{p}(w | t)) = \sum_{d, w} \frac{n_{dwt}}{\hat{n}_t} \ln \frac{n_{dwt}}{E_{dwt}}.$$

Статистика $G_t^2 = \sum_{d \in D} G_{dt}^2 = 2\hat{n}_t \text{KL}_t$ имеет асимптотически распределение χ_k^2 с числом степеней свободы $k = \sum_{d \in D} |W_{dt}| - |W| - |D| + 1$. В силу разреженности распределения $\hat{p}(d, w | t)$ вместо критерия хи-квадрат лучше применять перестановочный тест. Гипотеза условной независимости принимается для темы t , когда значение статистики G_t^2 меньше критического.

Выделение несогласованных тем. Статистические критерии позволяют находить «неудачные» темы, которые целесообразно разбивать на подтемы, непосредственно во время итераций EM-алгоритма. Темы можно ранжировать и сравнивать по значениям средней согласованности $B_t(\alpha)$ или статистики G_t^2 . Заметим, что сравнивать темы по значению дивергенции KL_t некорректно, так как только после умножения на «длину темы» \hat{n}_t получается величина $G_t^2 = 2\hat{n}_t \text{KL}_t$, имеющая (асимптотически) одинаковое распределение для всех тем.

§6.2 Критерии качества классификации документов

Оценивание качества тематической модели упрощается в тех случаях, когда она строится с целью классификации или поиска документов. Каждый документ описывается $|T|$ -мерным вектором тем $\theta_d = (p(t|d))_{t \in T}$. Качество модели определяется тем, насколько хорошо классифицируются документы, представленные этими векторами.

Пусть каждый документ $d \in D$ относится к классу $y_d \in Y$, алгоритм классификации $a: \mathbb{R}^{|T|} \rightarrow Y$ относит документ d к классу $a_d = a(\theta_d)$. В задачах информационного поиска и категоризации текстов качество классификации принято измерять в терминах точности и полноты [29].

Точность (precision) относительно класса $y \in Y$ определяется как доля правильно классифицированных документов среди всех документов, отнесённых алгоритмом a к классу y :

$$P_y(a) = \frac{\#\{d \in D: a_d = y_d = y\}}{\#\{d \in D: a_d = y\}}.$$

Полнота (recall) относительно класса $y \in Y$ определяется как доля правильно классифицированных документов среди всех документов класса y :

$$R_y(a) = \frac{\#\{d \in D: a_d = y_d = y\}}{\#\{d \in D: y_d = y\}}.$$

Чем больше значения точности и полноты, тем выше качество классификации.

В задачах информационного поиска обычно рассматривают два класса — документ либо «релевантен», либо «нерелевантен»; точность и полноту определяют только относительно класса релевантных документов.

Задачи категоризации, как правило, являются *многоклассовыми*, $|Y| \gg 2$. В таких случаях точность и полноту усредняют по всем классам.

В качестве агрегированного показателя, объединяющего точность P и полноту R , принято использовать F_1 -меру:

$$F_1 = \frac{2PR}{P + R}.$$

§6.3 Критерии качества тематического поиска

Описать идею разбиения каждого документа на части и поиска одних частей по другим. Качество поиска может измеряться с помощью Mean Average Precision. ToDo²⁷

Список литературы

- [1] Воронцов К. В., Потапенко А. А. Регуляризация, робастность и разреженность вероятностных тематических моделей // *Компьютерные исследования и моделирование*. — 2012. — Т. 4, № 4. — С. 693–706.
- [2] Лукашевич Н. В. Тезаурусы в задачах информационного поиска. — Издательство МГУ имени М. В. Ломоносова, 2011.

-
- [3] Маннинг К. Д., Рагхаван П., Шютце Х. Введение в информационный поиск. — Вильямс, 2011.
- [4] Павлов А. С., Добров Б. В. Метод обнаружения массово порожденных неестественных текстов на основе анализа тематической структуры // *Вычислительные методы и программирование: новые вычислительные технологии*. — 2011. — Т. 12. — С. 58–72.
- [5] Asuncion A., Welling M., Smyth P., Teh Y. W. On smoothing and inference for topic models // *Proceedings of the International Conference on Uncertainty in Artificial Intelligence*. — 2009.
- [6] Blei D. M., Ng A. Y., Jordan M. I. Latent Dirichlet allocation // *Journal of Machine Learning Research*. — 2003. — Vol. 3. — Pp. 993–1022.
- [7] Buntine W. L. Estimating likelihoods for topic models // *1st Asian Conference on Machine Learning: Advances in Machine Learning*. — 2009. — Pp. 51–64.
http://www.nicta.com.au/_data/assets/pdf_file/0019/20746/sdca-0202.pdf.
- [8] Celeux G., Chauveau D., Diebolt J. On stochastic versions of the em algorithm: Tech. Rep. RR-2514: INRIA, 1995.
- [9] Chemudugunta C., Smyth P., Steyvers M. Modeling general and specific aspects of documents with a probabilistic topic model // *Advances in Neural Information Processing Systems*. — MIT Press, 2006. — Vol. 19. — Pp. 241–248.
- [10] Cressie N., Read T. R. C. Multinomial goodness-of-fit tests // *Journal of the Royal Statistical Society, Series B*. — 1984. — Vol. 46, no. 3. — Pp. 440–464.
- [11] Daud A., Li J., Zhou L., Muhammad F. Knowledge discovery through directed probabilistic topic models: a survey // *Frontiers of Computer Science in China*. — 2010. — Vol. 4, no. 2. — Pp. 280–301.
<http://dx.doi.org/10.1007/s11704-009-0062-y>.
- [12] Dempster A. P., Laird N. M., Rubin D. B. Maximum likelihood from incomplete data via the EM algorithm // *J. of the Royal Statistical Society, Series B*. — 1977. — no. 34. — Pp. 1–38.
- [13] Eisenstein J., Ahmed A., Xing E. P. Sparse additive generative models of text // *ICML'11*. — 2011. — Pp. 1041–1048.
- [14] Feng Y., Lapata M. Topic models for image annotation and text illustration // *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*. — Association for Computational Linguistics, 2010. — Pp. 831–839.
- [15] Grün B., Hornik K. Topicmodels: An r package for fitting topic models // *Journal of Statistical Software*. — 2011. — Vol. 40, no. 13. — Pp. 1–30.
- [16] Hoffman M. D., Blei D. M., Bach F. R. Online learning for latent dirichlet allocation // *NIPS*. — Curran Associates, Inc., 2010. — Pp. 856–864.

-
- [17] Hofmann T. Probabilistic latent semantic indexing // Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval. — New York, NY, USA: ACM, 1999. — Pp. 50–57.
- [18] Kim S.-H., Choi H., Lee S. Estimate-based goodness-of-fit test for large sparse multinomial distributions // *Computational Statistics and Data Analysis*. — 2009. — Vol. 53, no. 4. — Pp. 1122 – 1131.
<http://www.sciencedirect.com/science/article/pii/S0167947308004817>.
- [19] Krestel R., Fankhauser P., Nejdl W. Latent dirichlet allocation for tag recommendation // Proceedings of the third ACM conference on Recommender systems. — ACM, 2009. — Pp. 61–68.
- [20] Li X.-X., Sun C.-B., Lu P., Wang X.-J., Zhong Y.-X. Simultaneous image classification and annotation based on probabilistic model // *The Journal of China Universities of Posts and Telecommunications*. — 2012. — Vol. 19, no. 2. — Pp. 107–115.
- [21] Lu Y., Mei Q., Zhai C. Investigating task performance of probabilistic topic models: an empirical study of PLSA and LDA // *Information Retrieval*. — 2011. — Vol. 14, no. 2. — Pp. 178–203.
- [22] Masada T., Kiyasu S., Miyahara S. Comparing lda with plsi as a dimensionality reduction method in document clustering // Proceedings of the 3rd International Conference on Large-scale knowledge resources: construction and application. — LKR’08. — Springer-Verlag, 2008. — Pp. 13–26.
- [23] Mimno D., Blei D. Bayesian checking for topic models // 11th Conference on Empirical Methods in Natural Language Processing. — Association for Computational Linguistics, 2011. — Pp. 227–237.
- [24] Minka T. P. Estimating a dirichlet distribution: Tech. rep.: 2000 (revised 2003, 2009, 2012).
- [25] Nigam K., McCallum A. K., Thrun S., Mitchell T. Text classification from labeled and unlabeled documents using em // *Machine Learning*. — 2000. — Vol. 39, no. 2-3. — Pp. 103–134.
- [26] Pecina P., Schlesinger P. Combining association measures for collocation extraction // Proceedings of the COLING/ACL on Main conference poster sessions. — Association for Computational Linguistics, 2006. — Pp. 651–658.
<http://http://dl.acm.org/citation.cfm?id=1273073.1273157>.
- [27] Read T., Cressie N. Goodness-of-Fit Statistics for Discrete Mutivariate Data. — Springer, New York, 1988.
- [28] Rubin T. N., Chambers A., Smyth P., Steyvers M. Statistical topic models for multi-label document classification // *Machine Learning*. — 2012. — Vol. 88, no. 1-2. — Pp. 157–208.

-
- [29] *Sebastiani F.* Machine learning in automated text categorization // *ACM Computing Surveys*. — 2002. — Vol. 34, no. 1. — Pp. 1–47.
- [30] *Steyvers M., Griffiths T.* Finding scientific topics // *Proceedings of the National Academy of Sciences*. — 2004. — Vol. 101, no. Suppl. 1. — Pp. 5228–5235.
- [31] *Taneichi N., Sekiya Y., Imai H.* Improvements of goodness-of-fit statistics for sparse multinomials based on normalizing transformations // *Annals of the Institute of Statistical Mathematics*. — 2003. — Vol. 55. — Pp. 831–848.
<http://dx.doi.org/10.1007/BF02523396>.
- [32] *Teh Y. W., Newman D., Welling M.* A collapsed variational bayesian inference algorithm for latent dirichlet allocation // *NIPS*. — 2006. — Pp. 1353–1360.
- [33] TextFlow: Towards better understanding of evolving topics in text. / W. Cui, S. Liu, L. Tan, C. Shi, Y. Song, Z. Gao, H. Qu, X. Tong // *IEEE transactions on visualization and computer graphics*. — 2011. — Vol. 17, no. 12. — Pp. 2412–2421.
- [34] *Vulić I., Smet W., Moens M.-F.* Cross-language information retrieval models based on latent topic models trained with document-aligned comparable corpora // *Information Retrieval*. — 2012. — Pp. 1–38.
- [35] *Wallach H.* Structured Topic Models for Language: Ph.D. thesis / Newnham College, University of Cambridge. — 2008.
- [36] *Wallach H., Mimno D., McCallum A.* Rethinking LDA: Why priors matter // *Advances in Neural Information Processing Systems 22* / Ed. by Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, A. Culotta. — 2009. — Pp. 1973–1981.
http://books.nips.cc/papers/files/nips22/NIPS2009_0929.pdf.
- [37] *Wallach H., Murray I., Salakhutdinov R., Mimno D.* Evaluation methods for topic models // *26th International Conference on Machine Learning, Montreal, Canada*. — 2009. — Pp. 1105–1112.
<http://www.cs.umass.edu/~mimno/papers/wallach09evaluation.pdf>.
- [38] *Wang Y.* Distributed Gibbs sampling of latent dirichlet allocation: The gritty details. — 2008.
- [39] *Wu Y., Ding Y., Wang X., Xu J.* A comparative study of topic models for topic clustering of chinese web news // *Computer Science and Information Technology (ICCSIT), 2010 3rd IEEE International Conference on*. — Vol. 5. — july 2010. — Pp. 236–240.
- [40] *Yeh J.-h., Wu M.-l.* Recommendation based on latent topics and social network analysis // *Proceedings of the 2010 Second International Conference on Computer Engineering and Applications*. — Vol. 1. — IEEE Computer Society, 2010. — Pp. 209–213.
- [41] *Yi X., Allan J.* A comparative study of utilizing topic models for information retrieval // *Advances in Information Retrieval*. — Springer Berlin Heidelberg, 2009. — Vol. 5478 of *Lecture Notes in Computer Science*. — Pp. 29–41.

-
- [42] *Zavitsanos E., Paliouras G., Vouros G. A.* Non-parametric estimation of topic hierarchies from texts with hierarchical Dirichlet processes // *Journal of Machine Learning Research*. — 2011. — Vol. 12. — Pp. 2749–2775.
<http://dl.acm.org/citation.cfm?id=1953048.2078193>.
- [43] *Zelterman D.* Goodness-of-fit tests for large sparse multinomial distributions // *Journal of the American Statistical Association*. — 1987. — Vol. 398, no. 82. — Pp. 624–629.
- [44] *Zhang J., Song Y., Zhang C., Liu S.* Evolutionary hierarchical Dirichlet processes for multiple correlated time-varying corpora // Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. — 2010. — Pp. 1079–1088.
- [45] *Zhang Z., Iria J., Brewster C., Ciravegna F.* A comparative evaluation of term recognition algorithms // Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC08). — 2008.
http://http://www.dcs.shef.ac.uk/~kiffer/papers/Zhang_LREC08.pdf.